Agristars

SR-P1-04194 NAS9-15466

A Joint Program for Agriculture and Resources Inventory Surveys Through Aerospace Remote Sensing

Supporting Research

December 1981

Technical Report

Multistage Classification of Multispectral Earth Observational Data: The Design Approach

by M.J. Muasher and D.A. Landgrebe

Purdue University Laboratory for Applications of Remote Sensing West Lafayette, Indiana 47907











General Disclaimer

One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.
- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.
- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.
- This document is paginated as submitted by the original source.
- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.

Produced by the NASA Center for Aerospace Information (CASI)

Agristars

"Made available under NASA sponsorship in the interest of early and wide dissemination of Earth Resources Survey Program Simation and without liability for any unit made thereot."

Supporting Research

E82 10213 SR-P1-04194, LARS 101481 S.NAS9-15466 CR-167508 A Joint Program for Agriculture and Resources Inventory Surveys Through Aerospace Remote Sensing

6. December 1981

Technical Report

Multistage Classification of Multispectral Earth Observational Data: The Design Approach

by M.J. Muasher and D.A. Landgrebe

Purdue University Laboratory for Applications of Remote Sensing West Lafayette, Indiana 47907

(E82-10213) MULTISTAGE CLASSIFICATION CF MULTISPECTRAL EARTH OBSERVATIONAL DATA: THE DESIGN APPROACH (Purdue Univ.) 181 p HC A09/MF A01 CSCL 02C

N82-23586

Unclas 00213









63/43



SR-P1-04194 NAS9-15466 LARS 101481

MULTISTAGE CLASSIFICATION OF MULTISPECTRAL EARTH OBSERVATIONAL DATA: THE DESIGN APPROACH

M.J. Muasher and D.A. Landgrebe

Purdue University Laboratory for Applications of Remote Sensing West Lafayette, Indiana 47907-05(1

December 1981

.

TABLE OF CONTENTS

Page

LIST	0F	T.	ABI	LES	5	•	•		•	•	•	•		•	•	•	•		•	٠	•	•	•	•	•	•	•	111
LIST	OF	F	IG	URI	ΞS	•	•		•	•	٠	•		•	•	•	•		•	•	•	•	•	•	•	•	•	ív
ABST	RA C'	Г	•	•	•	•	•		•	•	•	•		•	•	•	•	•	•	•	•	•	•	•	•	٠	•	vii
CHAP	TER	1	-	I	N T I	RO	DU	Cl	r I	01	1	•		•	•	•		1	•	•	•	•	•	•	•	•	•	1
	1.	1	Mu	Lt:	ĺs	ta	ge	(c 1	as	s	1 f	i	са	ti	Lο	n				•		•	•		•	•	1
	1.	2	Re	vie	ew	o	f	L:	1 t	eı	: a	tu	r	е					•		•	•				•		6
			1 '	2 ·	1	T τ	e f	n .	i n	0	Р	r o		e d	11 1	re		_			-							6
			+••	- • •	÷.	 -	- F	~		0			ี น	- t	-1 -		+ -		-	•	•	•	•	•	•	•	Ī	7
				~ ••	<u> </u>	ге И	1.1	101	L UI	a 1		ິດ	1	0 L 		ua L	4.4		3	•	•	•	•	•	•	•	•	12
		~	1	۷.	וכ	мu	TC	I	б С -	ag	;e	U	ι T	as	81	LI	16	: [8	•	•	•	•	•	•	•	•	10
	1.	-	Sui	nma	a r	y	of		Co	nt	:e	nt	s		٠	•		•	•	٠	•	٠	•	•	٠	•	•	19
CHAP	TER	2		P	AR.	ΑM	ΕТ	E	R	CC) N	SI	D	ER	A']	ΓI	10	١S	F	01	R							
				Α	M	UL	ΤI	S	ΤA	GE	2	BI	N	AR	Y	Т	RE	ΕE	C	CL/	A S S	SIF	TIE	R	•	•	•	20
	2.	1	The	e 1	Hu	gh	es	1	Рħ	er	10	me	n	on	L													20
	2	2	C 4 .		1 +	2 m			c	n f	10	00	n	<u>a</u> 1	1.	, ,	+ 1	0	n •		r h e	201	۰v					27
	2.	2 ว	5 II	.u.u.		a 11 -	60	1		د رز است		БV —		а т			-							•	•	•	•	20
	2.	<u>ر</u>	rea		ur:	e .	5e	10	ec	. L 1	10	n		•	•	•	•	•	•	•	•	•	•	•	•	•	•	29
	2.	4	S 11	mu.	La	ti	on		A I	go	r	1 t	: n	m	•	•	•	•	•	•	٠	•	•	•	٠	•	•	32
			2.4	4.	1	Ne	ed		Fο	r	A	S	1	mu	14	at	ic	o n	F	11	goı	: i t	:hπ	1	•	٠	٠	32
			2.4	4.:	2	St	a t	1	st	10	2.8	1	В	a c	kį	gr	٥ı	n	d	•	•	•	•	•	•	•	•	35
CHAP	TER	3	-	P	ER	FO	RM	(A)	N C	E	Е	SI	'I	MA	T) R	:	A	ΡĒ	R	oxi	I MA	TI	10	1	го		
				T	ΗĒ	Р	RO	B	A B	TI	.т	ΥY	,	0 F	' 1	ER	RC) R										37
						-														•	-	-	-	-				
	2	1	тъ	. ·		ka	14	h	~ ~	а	F			+ 4	~	-												37
	່ງ. າ	ר י	1 II D -	с. 		ле 		. 111	00	, u	- T - L	u 11 	i Ci	ь.l	. 01			•	•	•	•	•	•	•	•	•	•	60
	3.	2	re		οr -	ma 	nc	e	Ľ	. 8 1	[1	ma	ĽĽ	οr		•	. '	•	•	٠	•	•	•	•	•	•	•	42
			з.	2.	T	ſh	e	N	or	ma	11	A	S	s u	m	pτ	10	o n		٠	٠	٠	•	٠	•	٠	٠	43
			3.	2.	2	Тħ	e	M	o d	11 f	Εi	вç	l ć	Ga	m	na	Į	١s	sι	ımj	pti	ĹOI	1:					
						Fu	k u	in.	a g	,a	а	n ċ	1	Κr	: 1	l e	1	Ve	rs	s 1 (on	•		•	•	•	•	44
			3.	2.	3	Pr	op	0	sē	d		Mc	o d	1 f	:i(e d		1	go	or:	itl	hm	•			•		56
							-				-								-									

- Aller and the second second

đ

union and the desired of the state of the second state of the seco

i

· •

,

F	'age
CHAPTER 4 - EXPERIMENTAL RESULTS	64
4.1 Introduction	64
4.2 Experiments on Feature Selection	65
4.3 Experiments on the Hughes Phenomenon	74
4.4 Experiments Comparing Algorithm and	
Experimental Results	84
4.5 Experiments on a Binary Tree Classification	•
Procedure	104
	104
CHAPTER 5 - SUMMARY AND CONCLUSIONS	114
5 1 Summary of Pogults	114
5.1 Summary Of Results	116
J.2 Suggestions for Further Research	112
LIST OF REFERENCES] 1 7
APPENDICES	
Appendix A Generation of Normally Distributed	
Appendix A Generation of Normally Distributed	124
Appendix B. On The Drehability Density Functions	124
Appendix b on the probability bensity runctions a^2	107
$\begin{array}{c} \mathbf{OI} \mathbf{OI} \mathbf{And} \mathbf{OI} \mathbf{OI} $	121
Appendix C Derivation of the variance of	
o_1° and o_2°	134
Appendix D Classification Results Tables	141
Appendix E Computer Program Listings	150
Appendix F Description of Data Sets For	
Experiments	161
VITA	172

LIST OF TABLES

Table		Page
D.1	Classification Results of Aircraft, Simulated Data, Using 20 samples per class	142
D.2	Classification Results of Aircraft, Simulated Data, Using 13 samples per class	143
D.3	Classification Results of Aircraft, Real Data, Jsing 20 samples per class	144
D.4	Classification Results of Aircraft, Real Data, Using 13 samples per class	145
D.5	Classification Results of Landsat, Mul+itemporal, Simulated Data, Using 20 samples per class	146
D.6	Classification Results of Landsat, Multitemporal, Simulated Data, Using 13 samples per class	147
D.7	Classification Results of Landsat, Multitemporal, Real Data, Using 20 samples per class	148
D.8	Classification Results of Landsat, Multitemporal, Real Data, Using 13 samples per class.	149

•

.

LIST OF FIGURES

Figur	e	Page
1.1	An Example of a "Single-Stage" Algorithm In Classifying Multispectral Data	2
1.2	An Example of a "Multi-Stage" Algorithm In Classifying Multispectral Data	5
2.1	The Hughes Phenomenon	21
2.2	Explanation of the Hughes Phenomenon	2 2
2.3	Delineation of Optimal Subspace by Simultaneous Diagonalization	30
3.1	Probability Density Functions of h(X/w _i) and The Probability of Error	41
3.2	A Flowchart of Fukunaga and Krile's Algorithm	57
4.1	Classification Results of Data in Experiment 4.1 Using Three Feature Selection Techniques	68
4.2	Classification Results of Data in Experiment 4.2 Using Three Feature Selection Techniques	70
4.3	Classification Results of Data in Experiment 4.3 Using Three Feature Selection Techniques	72
4.4	Experimental Classification Results of Aircraft, Simulated Data Using Different Numbers of Training Samples	76
4.5	Experimental Classification Results of Aircraft, Real Data Using Different Numbers of Training Samples	79
4.6	Experimental Classification Results of Landsat, Multitemporal, Simulated Data Using Different Numbers of Training Samples.	81

ř.

1 v

F	i	g	u	r	e
---	---	---	---	---	---

4 J.A

4.7	Exp	еr	1	me	n :	ta	1	C	:1	88	36	i	fi	. c i	a t	1	or	1	Re	2 6	u	lt	5	c	>f		La	ın	d٤	5 a	t.			
	Mul	t i	t	en	p	01	a	1,		Re	ea	1	I	a)	ta		Ur	4	ng	3	D	i f	f	e	r e	n	t							
	Num	be	r	6	0	I	T	ra	11	n	l n	8	5	at	пp	1	ee	3.			•	•		•	•		•	•		•	•	٠		82
4.8	Cla	5 6	1	£i	c	a t	: i	or	1	Re	2 8	u	1 t	S	O	f	F	⁷ u	kι	ın	a	ga		ar	١đ	1	Κr	:1	16	e '	6			~ ~
	Еха	mΡ	1	e	R	ep	r	00	lu	C €	e d	•		•	•		•	•	1	•	•	•		•	٠		•	•		•	•	•		86
4.9	A F	10	w	ch	a	rt	. (o f		tł	ne]	Mo	di	l f	1	e d	I	A 1	lg	01	ri	t	ht	n		•	•	•	•	•	•		88
4.10	Cla	s 5	i	fi	. c :	a t	:1	on	1	Re	8	u	lt	s	o	£	A	1	r	r	a	ft	,	2	1	m	u 1	a	te	e d				
	Dat	а,	1	Vs	:11	n g	5	2 0)	S a	I m	p	le	S	Ρ	e	r	С	18	ıs	S	•		•	•		•	•	•	•	•	•		89
4.11	Cla	s s	1	fi	са	a t	: 1	០ព	1	Rε	2 S	u .	lt	S	o	£	A	.i	rc	r	a	E t:	,	S	i	m١	u 1	a	tε	e đ				
	Dat	a,	1	Us	:11	n g	5] 3	}	S a	m	р	le	5	Р	e	r	С	1 <i>a</i>	S	5	•	-	•	•		•	•	•	•	•	•		92
\$.12	Cla	ទម	i :	fi	c	a t	1	οn	1	Re	: 5	u]	l t	s	0	£	A	i	rc	r	ai	ft	,	F	le	a	L	Da	a t	a	,			
	Usi	ng		2 0) !	5 a	mj	p 1	е	S	p	e	r	C1	la	S	5.			,	•	•		•	•		•	•	•	,	•	•		93
4.13	Cla:	s s	1	f 1	C	a t	1	o n	1	Re	8	u I	l t	s	0	f	A	i	rc	r	a i	ft	,	F	le	a	L	Da	a t	a	,			
	Usi	n g	-	13	9	Sa	mj	p 1	e	S	P	eı	r	CI	la	S (s .		•		•	•		•	•		¢	٠	•	,	•	•		95
4.14	Cla	5 S	1:	fi	Cá	a t	i	o n	l	Re	s	u]	l t	s	0	f	L	a	n d	ls	a t	t,]	Mu	1	t	Lt	eı	mp	0	ra	1,	•	
	Sim	u1	at	te	d	D	a	ta	,	U	s	iı	ng	2	20	9	Sa	m	p 1	e	S	p	e	r	С	1:	15	8	•		•	•		97
4.15	Cla:	s s	1 :	fi	ca	a t	10	o n	L	Re	: 6	u]	ίt	s	0	f	L	a	n d	ls	a t	Ŀ,	1	Μu	1	ti	l t	e	np	0	ra	1,	,	
	Sim	u1	a 1	te	đ	D	at	ta	,	U	5	i r	۱g	1	. 3	5	Sa	m	р <u>1</u>	e	s	P	e	r	С	1 8	15	S	•		•	•		99
4.16	Cla	5 S	i 1	fi	ca	i t	i) n		Re	s	u]	lt	s	о	f	L	a	n d	ls	a t	Ŀ,	J	Mu	1	ti	l t	e	np	0	ra	1.	,	
	Rea	1	Da	a t	а,	,	U	s i	n	g	2	0	S	аn	Ъp	16	9 5	1	рe	r	(ci	a	s s	•		•	•	•		•	•	1	00
4.17	Cla:	55	i :	fi	ca	i t	10	o n		Re	5	u 1	Lt	s	0	f	L	a	n d	s	a t	Ξ.	1	Mu	1	ti	i t	er	np	0	ra	1.		
	Rea	1	Da	a t	а,	,	Us	si	n	g	1	3	S	аπ	ıp	16	2 5	1	p e	r	C	ci	a	5 5	•		•	•			•	•	1	02
4.18	Bina	ar	v	Т	re	e e	I	De	S	ig	n	ç	st	ru	١c	tι	ır	e	0	f	L	La	n	ds	а	t.								
	Mul:	ti	te	em	рc	r	a]	L,	1	Re	a	1	D	a t	a	,	U	s	in	g	1	3	•	Γr	a	i r	i i	n۶	3					
	Sam	p1	e 9	5	p e	2 r) N .		a:	5 5	,	4 •	11 :	th T-	1	N ı 4 -	ım	be	e r	S C	I	n	S : 1	i d	e	1 	10	de	s s	;			,	<u>م</u> د
	ING	T C	41	L]	11 8	5	NI	ш	שו	e r	(U I	-	ΤŢ	a	τĭ	11	11 (R	21	a 11	цр	T	ะ S		υ 5	ie	α.	•		•	•	Ŧ	00
4.19	Sin	g1	е	S	ta	a g	e	а	n	d	B	ir	na	ry	,	T 1	re	e	С	1	a s	ss	1	fi	C	a t	:1	or	ı					
	Resi	ul	t s	5 1 7	01	t 	La	a n	d	sa °	t	,	М -	u1	.t	11	te	m	p o	r	a 1	L,	Ī	Re	a	1	D	at	: a	,			4	^ 0
	USTI	ug		r D	- 2) a	шF	ιt	e	5	pe	e 1		υI	.а	58	÷.				•	٠		•	•						•		1	UΒ

Page

Figure	e	Page
4.20	Binary Tree Design Structure of Aircraft, Real Data, Using 13 Samples per Class	110
4.21	Single-Stage and Binary Tree Classification Results of Aircraft, Real Data, Using 13 Training Samples per Class	111

v1

•

ABSTRACT

One of the main problems in a multistage decision tree procedure is predicting the optimal features to be used at every node. An algorithm is proposed which predicts the optima' lectures at every node in a binary tree procedure. The algorithm estimates the probability of error by approximating the area under the likelihood ratio function for two classes, and taking into account the number of training samples used in estimating each of these two classes. Some results on feature selection techniques, particularly in the presence of a very limited set of training samples are presented. Results comparing probabilities of error predicted by the proposed algorithm as a function of dimensionality as compared to experimental observations are shown for aircraft Results are obtained for both real and and Landsat data. simulated data. Finally, two binary tree examples which use the algorithm are presented to illustrate the usefulness of the procedure.

2 63.22 may property

CHAPTER 1

INTRODUCTION

1.1 Multistage Classification

A number of different types of classifiers are now in routine use in remote sensing. Most of these classification algorithms, using pattern recognition techniques, can be regarded as "single-stage" classifiers, where an "unknown" pattern is tested against all classes using one feature subset, and then the pattern is assigned to one of the present classes in a single-stage decision procedure. An example of such a procedure is shown in Figure 1.1.

In recent years, as classification of multispectral data has found a larger number of users and a wider range of applications, the need has been felt for alternate, more powerful techniques than the conventional classifiers, through the use of which more information could be extracted more accurately and/or efficiently from the scene. Some of the reasons that have warranted this need include:

 The need to extract more detailed information from data. The opportunity to do so results from the



Figure 1.1 An Example of a "Single-Stage" Algorithm In Classifying Multispectral Data.

emergence of more complex data sets. The growing use of multitype data bases containing Landsat data with a variety of other quantitative geodata together with the anticipated launching of more sophisticated sensors such as the Thematic Mapper result in the opportunity to extract considerably more information from the data.

- 2. The broadening of the range of applications. As pattern recognition methods have developed, they have found a larger number of users with a wider range of applications. The feedback from these different and versatile uses has indicated problems and needs not initially present.
- 3. The ever present need for improved classification accuracy. There are some applications for which conventional classifiers have proved to be marginal at best. Some of these are listed in Swain et al. (1) and include multi-image analysis and the use of mixed feature types.
- 4. The need for improved processing efficiency. The conventional, single-stage, classifiers use only one parricular feature subset and are somewhat inefficient, as they must compare an unknown pattern against all possible classes before assigning that pattern to a particular class.

Because of these and other factors, there has been some research in recent years directed towards developing multistage classifiers, whereby the decision procedures go through several stages before finally assigning a pattern to a class. An example of such a procedure is shown in Figure 1.2.

The purpose of this research is to develop a layered decision algorithm that can increase the accuracy and efficiency over the conventional single-stage classification approach. Developing such an algorithm requires, among other things, a careful look at some parameters that are crucial to any successful attempt at tackling such a complex problem. In particular, three areas have to be investigated:

- The development of an adequate training procedure to define an initial set of spectral classes with their respective statistics;
- 2. The investigation of various error estimators and the development of an adequate performance estimator that can reasonably predict the accuracy or any trends in performance;
- 3. The development of an algorithm to build a binary tree making use of the above-mentioned methods.

Martin Land & Constant Street Street



Figure 1.2 An Example of a "Multi-Stage" Algorithm In Classifying Multispectral Data. Of these three areas, the most important problem is believed to be the development of an accurate error estimator, especially in the presence of what has come to be known as the Hughes phenomenon (elaborated upon later in the review of literature). Predicting the conditions under which the Hughes phenomenon occurs provides the key to the solution of the problem. Therefore, a considerable portion of the research has been directed towards trying to understand and predict the impact of this phenomenon.

1.2 Review of Literature

1.2.1 Training Procedure

Several training methods have been suggested in the literature. We will not attempt to list all of them, but rather will give a background of some of the methods reviewed and used in this work.

The training process is the procedure whereby labeled samples are selected and used to compute class statistics which in turn are used to classify unlabeled (i.e., "unknown") samples.

Several parameter estimation methods (training methods) have appeared in the literature. Sample-partitioning methods, the leaving-one-out method, clustering are but a few. See, for example, Fukunaga (2) and Duda and Hart (3).

For remote sensing purposes, clustering has been widely developing training statistics. Two basic used in approaches have been: a supervised clustering approach, in which the analyst selects areas of known cover types , each set of areas belonging to one cover type is clustered sepaand then the statistics for these areas are then rately. obtained with the aid of a computer; and the non-supervised clustering approach, in which the entire training area is subdivided into clusters by the clustering algorithm and each cluster is then identified by the analyst and given a specific label. The statistics of each cluster corresponding to a cover type or a subclass of a cover type are then calculated. Fleming et al. (4,5) investigated several clustering approaches and their effect on classification accu-Among the approaches they used were non-supervised racy. clustering, supervised clustering, modified clustering, mono- (aggregate) cluster blocks, and multi- (class-conditional) cluster blocks.

1.2.2 Performance Estimators

A key factor in the design of a layered decision algorithm is the ability to predict how the algorithm will perform in terms of accuracy at every node. While optimizing the performance at every node does not necessarily produce a globally optimal tree, it is still a very important and useful step in the design.

Several performance (or error) estimators have appeared in the literature. Again, we will not attempt here to exhaustively list all the contributions made, but rather will give an idea of how the research in this area has progressed.

Performance estimators can be divided into two main categories:

Performance functions which have some sort of direct relationship with the probability of error. Examples are Parzen estimators (see (2)), the k-nearest neighbor error estimator (see (6)). More recently, Mobasseri et al. (7) published an error estimator that computes the minimum probability of error through use of a combined analytical and numerical integration over a sequence of simplifying transformations of the feature space. The results have been shown to be similar to those obtained by conventional techniques. However, the algorithm becomes computationally too inefficient to use as the number of classes and/or features increases. Moore, Whitsitt and Landgrebe (8) (see also Whitsitt and Landgrebe (9)) developed a stratified posterior estimator which, like Mobasseri's, depends only on a given set of statistics. This was later used by Wiersma (10) and both estimators (Mobasseri's and Whitsitt's) were compared in (11) and found to give similar results, with Whitsitt's algorithm being faster in some cases. The former procedure

uses a "deterministic" grid to sample the feature space, while the latter uses an internally generated random data base and assigns the feature vector to the appropriate class via the maximum a posteriori principle. Both procedures assume normal class conditional statistics.

<u>Separability measures</u>, <u>most of which have only a sub-</u> <u>tle</u>, <u>indirect</u>, <u>and often unknown</u>, <u>relationship to the proba-</u> <u>bility of error</u>. Various separability measures have been in common use in remote sensing applications. Among these are: Divergence (12), Transformed Divergence (13), Jeffreys-Matusita distance (14,15), Bhattacharyya distance (16) and the Mahalanobis distance (17). (See list in (24).)

Several works have been reported comparing different separability measures and their effects on performance. (See (9,13,18,19,62).)

There are two problems with most of the above separability measures applied to remote sensing applications: (1) ambiguity and (2) linearity in pairwise error. The term ambiguity implies here that there dress not exist a one-toone relationship between the value of the measure and the probability of error. Linearity means that equal incremental changes in the measure imply equal changes in the probability of error, over the whole range. Whitsitt (9) developed a distance measure D erf = erf ($\sqrt{2B}$) where B is the Bhattacharyya distance and erf(.) is the gaussian error

Another key factor in the process of error estimation is the choice of feature subsets. The problems here are twofold:

1. As the number of features becomes large, it becomes desirable to choose a subset of these features that can adequately predict the accuracy. This selection process also can become expensive if one must search through all possible combinations of the feature set. It is desirable, therefore, to have a priori knowledge of the importance of each feature in relation to the probability of error. The Karhunen-Loeve expansion (attributed to Karhunen (20), and Loeve (21)) in pattern recognition literature has historically been used as a feature selection technique. It has the advantage of producing uncorrelated features (in theory, but the features are actually approximately uncorrelated in a practical K-L transformation). In addition, it imposes an ordering on the features in terms of importance in a representation error sense. As a result, first feature is "likely" to be more important than the second in calculating the probability of error, and so on. More recently, Oja and

Karhunen (22,23) published two papers on the construction of K-L expansions for pattern recognition purposes that do not require the computation of any covariance matrices.

2. The probability of error is not necessarily monotonically decreasing as the number of features increases. This is due to a peculiar phenomenon that has come to be known as the Hughes phenomenon. Hughes (25) found that with a fixed and finite training pattern sample, recognition accuracy can first increase as the number of measurements on a pattern increases, but decay with measurement complexity higher than some optimum value. He also reported that for unlimited training data, this does not occur and the recognition accuracy reaches an optimum only at infinite measurement dimensionality. According to Hughes, if insufficient sample data are available to estimate the pattern probabilities accurately, then a Bayes recognizer is not necessarily optimal. Many papers have since been published on this phenomenon, confirming it or trying to explain why it occurs (see (26-42)).Thus, it appears that a successful design should predict when and if such phenomena occur.

Constant to contractor and the second

1.2.3 Multistage Classifiers

TRANSPORT FOR ALL TO COMP

In recent years, some work has appeared in the literature aimed at developing multistage classification algorithms. There is much yet to be learned about such algorithms, and no work has been reported claiming optimality (or even close to optimality) of results.

In general, earlier work can be grouped into two main categories:

Sequential classification methods. These can be found in several papers and books (see, for example, (43-45)). Basically, the method consists of observations made on feature measurements, one at a time. After an observation is made, the classifier either reaches a final decision and the process is terminated, or it makes another observation until a final decision is rea hed.

<u>Hierarchical classification methods</u>. These are subdivided into two categories:

1. Hierarchical clustering methods. Examples of such work are found in Fukunaga (2), Dubes and Jain (46), who present a semi-tutorial review of the state of the art in cluster validity, and Lukasova (47). In general, hierarchical clustering is designed to generate a classification tree. The "root" node of the tree represents a collection of samples (either a training data set or the entire sample

set) and each terminal node represents either an individual sample or a group of samples belonging to some class within the set of classes in the data set. The method attempts to divide the set of samples in each node into disjoint subsets which form new nodes. Defined as such, the method is often nonparametric and depends heavily on the ability of the algorithm to find meaningful divisions of samples that correspond at terminal nodes with reaningful classes.

2. Decision trees and criterion functions. Most of the work done in multistage algorithms belongs to this category. Often, a decision tree is built using an optimization or criterion function that dictates the structure of the tree. It is this kind of approach that will be of greatest concern in this research.

Hierarchical methods differ from sequential methods in certain important respects. While in sequential schemes any class can be accepted at any stage of the measurement process, in hierarchical schemes certain classes are excluded from consideration at each scage. Also, sequential methods impose a linear ordering on the features. In hierarchical methods, features used along one decision path can be different from those used along another path.

In 1971, Nadler (48) tried to calculate error rates in a hierarchical decision structure under assumptions of statistical independence among the members of the hierarchy.

where the construction was as as we

Even under such assumptions, the results assume "small" probabilities of errors at any level.

Several heuristic methods of constructing tree designs have been proposed in the literature. Some studies were done using optimization methods to automate the classifier design procedure, but the assumptions made were often too restrictive. Meisel and Michalopoulos (49) in 1973 presented a two-stage partitioning algorithm for the design of an optimal binary tree. In the first stage, a suboptimal sufficient partition is obtained. The second stage optimizes the result of the first stage through a dynamic programming approach. The method allows only for linear discriminant functions to partition the space, certainly a suboptimal and tco restrictive condition.

In 1974, Wu et al. (50) reported on a decision tree approach with direct application to multispectral data analysis. Several design procedures were proposed (one of which is manual). with special emphasis on a heuristic, machine-implemented approach. The optimality criterion used is a weighted sum of computation cost and accuracy. Results were presented which showed superiority in efficiency (but infrequently in accuracy) over the conventional classifier. The criterion function used, as it cannot predict beforehand the structure of the tree below that node, assumes all the nodes below the node under consideration are terminal nodes,

and hence is necessarily suboptimal. Later papers have appeared that have pointed to applications using this particular classifier (51,52).

In 1976, You and Fu (53) presented a linear binary tree classifier that uses linear discriminant functions at decision stages with an application to multispectral remotely sensed data. The procedure includes a grouping algorithm, a separability measure, and an error minimization procedure using the Fletcher-Powell algorithm (54). Again, the procedure is certainly suboptimal because of the assumption of linearity. Results reported, though, show that this classifier is much faster and more accurate than the maximum likelihood classifier with the same number of features. This is due to the fact that the procedure uses different feature subsets (with a restriction on their number) at each node, compared with only one feature subset used in the one-stage maximum likelihood classifier.

Kulkarni and Kanal (55) used dynamic programming and branch-and-bound methodologies in the design of hierarchical classifiers. The criterion of optimality they used is a weighted sum of the probability of error and the average measurement cost incurred in classifying a random sample. The design assumes that the features used at the nodes are statistically independent and that the decision at each node is a function of only that particular feature observation,

the design using only one best feature at each tree node. Further, the design of the optimal tree assumes a very low error rate for the tree, a very restrictive assumption since in many cases a high error rate is specifically the reason why a layered classifier was selected, i.e., to improve the accuracy. Although the authors presented some methods to reduce the complexity of their design algorithms, the examples they used involve only a small number of classes and features.

In 1977, Parkih (56) compared several classification techniques of clouds, including hierarchical design. How-~ver, his paper offers no new insights or major results that would help improve the state of the art.

Also in 1977, Sethi and Chatterjee (57) developed an algorithm for the design of an efficient decision tree with application to pattern recognition problems involving discrete variables. A criterion function was defined to estimate the minimum expected cost of a tree in terms of the weights of its terminal nodes and costs of the measurements, which then was used to establish the search procedure for the efficient decision tree. The concept of prime events was used to obtain the number of nodes and the corresponding weights in the design sample. No optimality claim was made, but the procedure was found to lead to the optimal tree in most of the cases. The procedure uses only one feature at

Sector Contractor State Contractor

every node, and its applicability to remotely sensed multispectral data is very doubtful.

In 1978, Breiman (58) presented a procedure for building a binary classification tree. He used a criterion function that is only a function of the parent node and the two descendent nodes. He used one best feature at every node. He also reported on another regression algorithm developed at Survey Research Center, University of Michigan (59), in which the criterion function tries to reduce the variances of the two descendent nodes as much as possible from the variance of the parent node.

Rounds (60) in 1979 developed a binary decision tree algorithm, but again one feature is selected at every node. The approach is a nonparametric one, based on the Kolmogorese Smirnov criterion.

Dattatreya and Sarma (61) in 1981 presented a multistage binary tree "minimum-cost" classifier, when general cost functions are associated with the tasks of feature measurements. The optimization of the binary tree is carried out using dynamic programming. However, one feature is only selected at every node.

In summary, most of the work done with multistage classifier: often imposed too restrictive assumptions or conditions, such as using one feature only at each node. or hav-

ing a linear discriminant function. Moreover, very few results have been reported on situations where the Hughes phenomenon occurs, namely, working with a limited set of training samples.

The major contributions of this research are then:

- 1. The development of some theroretical results that clearly show the dependence of the accuracy of the estimated statistics of the classes under consideration on the number of training samples used to estimate the statistics of those classes, as well as on the number of features used.
- 2. The development of an error estimator which is particularly useful when the number of training samples is limited, and which is suited for a binary tree classification procedure. This estimator, which allows the selection of a "near optimal" feature subset at every node, has no restrictions on the number of features that can be used at any node.
- 3. The incorporation of the above error estimator in a binary tree procedure, showing the usefulness of such a procedure in predicting the optimal features that lead to the best accuracy that can be attained given a fixed sct of training samples.

1.3 Summary of Contents

In chapter 2, some parameter considerations for a multistage binary tree classifier are addressed in detail. The Hughes phenomenon is elaborated upon, and a technique known as "sumultaneous diagonalization" is introduced. Feature selection techniques are also treated. A data simulation algorithm that is repeatedly used in the research is also treated.

In chapter 3, an approximation algorithm to the probability of error is proposed that takes into account the Hughes phenomenon.

Chapter 4 presents experimental results on real and simulated data.

Finally, chapter 5 summarizes conclusions about the study. Some analytical details, together with computer listings and training data are placed in appendices.

CHAPTER 2

PARAMETER CONSIDERATIONS

FOR

A MULTISTAGE BINARY TREE CLASSIFIER

2.1 The Hughes Phenomenon

One of the major needs for a decision tree classifier originates from a dimensionality problem often referred to as the Hughes Phenomenon (25). A considerable portion of this research is directed towards understanding the Hughes phenomenon. Figure 2.1 illustrates the phenomenon conceptually. In the presence of a limited training sample size, the mean recognition accuracy as a function of the measurement complexity (number of features for our purposes) exhibits a peaking effect. Contrary to intuition, the mean accuracy does not always increase with additional measurements. Further, peaking of the curve shifts up and to the right as the number of samples increases, disappearing in the case of an infinite number of training samples (complete knowledge of the underlying distributions).

Figure 2.2 suggests a concept for one possible explanation of this phenomenon. Figure 2.2a shows a hypothetical



Figure 2.1 The Hughes Phenomenon.





2.2c

ł



graph of class separability plotted vs. dimensionality. As dimensionality increases, so does class separability (a nondecreasing function of dimensionality) until it saturates, and any further increase in dimensionality does not have a significant effect on class separability. But this is not the only effect on the mean accuracy. With the presence of a fixed, limited training sample size, any increase in dimensionality necessarily results on the average in a degradation in the accuracy of statistics estimation of the class distributions. Thus, conceptually, one should expect a curve similar to that of Figure 2.2b.. Further, as the number of samples increases, the curve should shift to the right, i.e., for any given dimensionality, the larger sample size should provide a better estimate of the true distributions. Assuming these two effects are the dominant effects on accuracy, adding the two effects results in Figure 2.2c, a curve similar to Figure 2.1. Based upon this concept of the phenomenon, the solution to the problem lies in being able to predict quantitatively how the number of samples present affects the accuracy of the estimated statistics . Especially in remote sensing applications of pattern recognition methods, training samples are limited as ground truth is often not present or difficult to get. Thus, the importance of the Hughes phenomenon becomes evident, as well as the validity of this conceptual explanation of it.

The Hughes phenomenon was studied by many researchers. (See (26-42)). Hughes (25), who was one of the earliest to introduce it and treat it in some detail, tried to explain it from a nonparametric point of view. The explanation given by Wacker and Landgrebe (62) is of another nonparametric case, where the Euclidean distance measure is used for discrimination among classes.

Several researchers (28-34) tried to study the effect of limited training sample size and independence of measurements on the recognition accuracy.

In 1979, Trunk (38) provided a simple example in which he showed theoretically that the probability of error approaches zero as the dimensionality increases and all the parameters are known in a two-class problem, but it approaches one-half as the dimensionality increases and the parameters are estimated.

In remote sensing applications, where maximum likelihood classifiers are frequently used, and where the assumption of class-conditional multivariate normally distributed data is invoked, not much work concerning the dimensionality problem has been reported yet. Wacker and El-Sheikh (40-42) presented some papers dealing with dimensionality problems for two-class Gaussian problems. Their results again show a Hughes phenomenon occuring with finite training data. It then follows that any error estimator in a multistage classification algorithm that can claim some optimality in results from an accuracy point of view, should be able to predict when/if a peaking occurs in the curve mentioned earlier. It is this key problem that this research is attempting to solve, i.e. the development of an error estimator that can accurately predict the Hughes phenomenon.

Working with multispectral data, one almost always has to work with multiple feature measurements and multiple classes. In this research, we propose a binary tree multistage classifier. This means that any node in the tree is either a terminal node or is further subdivided into two nodes (with statistics corresponding to two classes).

The advantages of a binary tree procedure are the following:

- 1. Working with two classes allows a theoretical understanding of the problem. Many pattern recognition results that apply to two-class problems fail to do so in multi-class ones. This is particularly true in the "simultaneous diagonalization" technique that will be introduced snortly.
- 2. Most feature selection algorithms used in pattern recognition applications generally, and in remote
sensing applications specifically, are optimal only when applied to two-class problems. For multiclass problems, a separability criterion is averaged over pairs of classes and thus is optimal only in an average sense. Working with a binary tree, then, should provide us with both convenience and accuracy.

Working with multiple features, several properties are desired in these features which will make further analysis easier:

<u>Uncoupled</u> (<u>Independent</u>) <u>Features</u>. Uncoupling of features from one another simplifies analysis a great deal as it permits evaluating the effect of each feature separately from other features.

Ordered Features. If the features can be ordered, or at least approximately so, in terms of their effect on the probability of error, then the process of feature selection would be made easier.

Optimal Separability. The features should be optimal with respect to the probability of error for two distributions at hand. Putting it in different words, the feature subset should be tailored to the separability of the two distributions. To this end, a technique known as a "simultaneous diagonalization" (63,64) is discussed in the next section.

2.2 Simultaneous Diagonalization: Theory

Let $\hat{\Sigma}_1$ and $\hat{\Sigma}_2$ be the estimated covariance matrices for classes 1 and 2, respectively. We seek a transformation matrix A such that

$$A \hat{E}_1 A = I$$
 $A \hat{E}_2 A = \Lambda$ (2.1)

where I is the identity matrix and Λ is a diagonal matrix.

This transformation would uncouple the features, while not affecting the probability of error because the latter is invariant under linear transformations. We proceed to find such a transformation as follows. (For more details, see (2), pp. 31-35.)

Let Θ and Φ be the eigenvalue and eigenvector matrices of $\hat{\Sigma}_1$, respectively; then

$$\Theta^{-\frac{1}{2}} \Phi^{T} \hat{\Sigma}_{1} \Phi \Theta^{-\frac{1}{2}} = \mathbf{I} \qquad (\Phi^{T} \hat{\Sigma}_{1} \Phi = \Theta) \qquad (2.2)$$

$$\Theta^{-\frac{1}{2}} \Phi^{T} \hat{\Sigma}_{2} \Phi \Theta^{-\frac{1}{2}} = K \quad \text{K is a general matrix} \qquad (2.3)$$

Next, we desire to diagonalize K. To find eigenvalues of K, it is necessary to solve the equation

$$|K - \lambda I| = 0 \tag{2.4}$$

Replacing K and I in (2.4) by (2.2) and (2.3), we get

$$\left[e^{-\frac{1}{2}} \phi^{T} \hat{\Sigma}_{2} \phi e^{-\frac{1}{2}} - \lambda e^{-\frac{1}{2}} \phi^{T} \hat{\Sigma}_{1} \phi e^{-\frac{1}{2}} \right] = 0$$
 (2.5)

 $\mathbf{0r}$

$$\left[\Theta^{-\frac{1}{2}} \phi^{\mathrm{T}} \right] \left[\hat{\Sigma}_{2} - \lambda \hat{\Sigma}_{1} \right] \phi = 0 \qquad (2.6)$$

Since $o^{-1_2} \phi^T$ is nonsingular, it follows that

$$\left| \hat{\Sigma}_{2} - \lambda \hat{\Sigma}_{1} \right| = 0$$
(2.7)

or,

$$\left| \hat{\Sigma}_{1}^{-1} \hat{\Sigma}_{2} - \lambda \mathbf{I} \right| = 0$$
(2.8)

So, only the eigenvalue and eigenvector matrices of $\hat{\Sigma}_1^{-1}\hat{\Sigma}_2$ need be calculated.

The eigenvalue matrix is then Λ , and the transpose of the eigenvector matrix, A^{T} , serves as the transformation matrix.

The idea behind simultaneous diagonalization is to transform the original features into a new space where the features are independent and then choose a subset of these features in the new space which is optimal with respect to the probability of error. This is illustrated in Figure 2.3.

2.3 Feature Selection

Before proceeding to discuss the approximation algorithms to estimate the probability of error, we digress briefly to discuss how the features are ordered.

The literature offers many studies made on comparing different separability measures and their effectiveness in choosing the best feature subset (see (9,13,18,62,55)). It appears that the Bhattacharyya distance is one of the most suitable separability measures for distinguishing between classes. Thus, it will be used as a basis for feature selection. The fact that the features are independent allows us to determine the effect of each feature on the probability of error separately.









The Bhattacharyya distance for two normal distributions can be expressed as follows:

$$B = \frac{1}{8} (\hat{M}_1 - \hat{M}_2)^T (\frac{\hat{\Sigma}_1 + \hat{\Sigma}_2}{2})^{-1} (\hat{M}_1 - \hat{M}_2) + \frac{1}{2} \ln \left| \frac{\frac{1}{2} (\hat{\Sigma}_1 + \hat{\Sigma}_2)}{\hat{\Sigma}_1 + \hat{\Sigma}_2} \right|$$

After the simultaneous diagonalization transformation, however, B can be expressed as:

$$B = \sum_{i=1}^{p} \left[\frac{1}{4} \frac{\left(\frac{d_{1i} - d_{2i}}{\lambda_{i} + 1}\right)^{2}}{1} + \frac{1}{2} \ln \left(\frac{1}{2} \left(\frac{1}{\lambda_{i}^{l_{2}}} + \lambda_{i}^{l_{2}}\right) \right) \right]$$
(2.10)

where d_{ij} is the jth element of the transformed class-conditional mean: $D_i = \stackrel{T}{A} \stackrel{\gamma}{M}_i$; and λ_i is the ith diagonal element of Λ .

Thus, it is clear that for every feature i, B can be calculated separately. The feature with the largest B is the best feature, the one with the second largest is the second best, and so on. Also, the two best are the best two, and so on. 2.4 Simulation Algorithm

2.4.1 Need For A Simulation Algorithm

For remote sensing data analysis, several assumptions are commonly made. These assumptions are usually that the data are class-conditionally distributed multivariate normal and that the data used to train the classifier are representative of the area of interest. This second assumption actually has several parts. The assumption is made that in the process of training, all classes present in the scene are found, and all spectral subclasses of each class are also represented in the training data. Furthermore, the parameters of the distribution of each subclass are also assumed to be known from the training data. Each pixel is assumed to come from one of the training classes, and also is assumed to be entirely of one cover type.

In actual practice, these assumptions are not met. The number of spectral classes in the area is not known and clustering or some other method is used to determine the number of subclasses, in addition to estimating the statistics of those subclasses. Some of these methods also lead to non-normal subclasses. In particular, the clustering algorithm available through LARSYS truncates the tails of the subclass distributions and so leads to non-normal distributions.

There are also questions relating to a single picture element. A single pixel in Landsat data covers an area approximately 80 meters by 50 meters. More than one cover type may be present in this area and result in a "mixture pixel" observation. It is not clear how the distribution of the spectral response of mixture pixels can be related to the distribution of the spectral response of "pure pixels".

There has been much speculation in the remote sensing community as to the effect of the non-satisfaction of the basic assumptions. Whenever new algorithms are brought forth, the old question is raised again, indicating that there is insufficient understanding of the interaction of the real attributes of the data and the theory of the algorithms. At times it is not clear whether a particular result is due to aspects of the algorithm or to the extent the data set deviates from the assumptions.

In testing new algorithms, deviations from the assumptions may obscure the action of the new process. One way to clarify the situation is to apply the algorithm first to a data set satisfying the assumptions.

Such a data set could be obtained artificially, through simulation. The analyst could then know: how many classes exist in the data; the true distributions of the classes, including normality if desired; the observations could really be independent; and no pixel would be a "mixture

المالية المائية المرجب والمرجب المرجب

pixel". New algorithms could be studied on such a data set with the knowledge that any "strange" effects are indeed algorithm rather than data problems.

In many cases where simulated data have been used in the past, the data were too artificial, in the sense that all aspects of the image were controlled, removing the natural variation in object size, position, and relationship which occur in real data. This limited the use of the simulated data sets in testing new algorithms.

The natural spatial information occuring in multispectral data could be retained in a simulated image by spatially basing the simulation on a classification. It would be even better to base the simulated data on a digitized "ground truth" map if the spectral characteristics of the cover types were known. By basing the simulation on a classification, the number of classes, their exact distributions. and the class of each pixel in the area are known. If the classification was sufficiently accurate, then the spatial information held in the classification map will be close to the actual cover type map and actual spatial content of the original data. For each pixel in the area. a random vector distributed according to the pixel's class statistics could be generated. This becomes the simulated data vector.

This simulated method was reported in LARS Technical Report 070980 (66), and the program will be used for testing the error estimator developed.

2.4.2 Statistical Background

From the classification chosen as a basis for the simulation, the following are known: the number of classes K, the set of classes $(\omega_i, i=1,\ldots,K)$, the class distributions $(f(\omega_i),i=1,\ldots,K)$, their means and covariances (ν_i and Σ_i , $i=1,\ldots,K$), the number of channels p, and the class of every pixel in the scene.

From classical statistics:

(1) Let X:px1, A:pxp, and b:px1.

If $X \sim N$ (0,I_p), then $Y = AX + b \sim N$ (b, $AI_pA^T = AA^T$) (where I is the identity matrix having dimensionality p).

(2) Let Σ be a symmetric, positive definite matrix. Then there exists A, such that $AA^{T} = \Sigma$ (A is denoted $\Sigma^{\frac{1}{2}}$)

To simulate a pixel which was a member of class i in the base classification, $N(0, I_p)$ (the random vector for each pixel is independent of other vectors) is generated. (See Appendix A.) Next $Y = \Sigma_i^{l_2} X + \mu_i$ is calculated; it is then a random vector from the population $N(\mu_i, \Sigma_i)$. This process is repeated for each pixel of the base classification and the random vectors thus generated are stored appropriately, i.e., so as to correspond to their simulated spatial location.

The program requires as an input a classification map stored on a results tape. The results tape has the class statistics for p-dimensions also stored on it. The program then, uses the results map and the stored statistics to generate a p-dimensional data set, which is stored on a user specified output tape in LARSYS format.

Appendix A provides a mathematical derivation related to the generation of normally distributed samples. Appendix E provides a Fortran program listing for the simulation program.

With all the proliminaries discussed, we are now ready to begin our discussion of the error estimator algorithm.

CHAPTER 3

PERFORMANCE ESTIMATOR:

APPROXIMATION TO THE PROBABILITY OF ERROR

3.1 The Likelihood Function

As mentioned earlier, our goal is to develop a performance estimator that can predict where the peak in the Hughes curve occurs. Some of the most serious difficulties facing researchers in trying to estimate the probability of error in multidimensional analysis are:

- 1. The need to carry out a multiple integration on the multivariate probability density function. Most often, this integration is almost impossible to carry out analytically, and numerical integration that is often costly has to be perfomed.
- 2. The measurement features are often correlated, making it difficult to assess the importance of each feature separately on the probability of error.
- 3. In most of the cases, one has to deal with multiclass problems (greater than 2) which further complicates multivariate probability density functions.

It would be much easier, therefore, if one could work with a function that is one-dimensional but carries all the information present. Fortunately, since we are looking at two classes at a time in a binary tree procedure, such a function does exist, and is called the likelihood function (minus the log of the likelihood ratio). See, for example, (66).

The likelihood function, denoted h(X), is given by:

$$h(X) = -\ln p(X/w_1) / p(X/w_2)$$
(3.1)

where

$$p(X/w_i)$$
 is the probability density function of X given w_i .

In remote sensing applications, the assumption of multivariate class-conditional normal distributions is almost always invoked, and will be consistently used in this work.

Using this assumption, $p(X/w_i)$ becomes:

$$p(X/w_{i}) = \frac{1}{(2\pi)^{p/2} \left| \Sigma_{i} \right|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (X^{T} - M_{i}^{T}) \Sigma_{i}^{-1} (X - M_{i}) \right) (3.2)$$

where M_i is the mean vector of class i.

 Σ_{i} is the covariance matrix of class i. p is the number of dimensions. In practice, M_i and Σ_i are estimated from training statistics and are replaced by \hat{M}_i and $\hat{\Sigma}_i$.

The Bayes decision rule for minimum error may be written as follows:

$$P(w_1/X) \stackrel{\geq}{\leftarrow} P(w_2/X) \rightarrow X \stackrel{\leftarrow}{\leftarrow} \begin{cases} w_1 \\ w_2 \end{cases}$$
(3.3)

The a posteriori probabilities $P(w_i/X)$ may be calculated from the a priori probabilities $P(w_i)$ and the conditional density functions $p(X/w_i)$ using Bayes theorem, i.e.

$$P(w_i/X) = p(X/w_i) P(w_i) / p(X)$$
 (3.4)

Since p(X) is common to both sides of the inequality of (3.3), the decision rule can be expressed as:

$$p(X/w_1) P(w_1) \stackrel{\geq}{\leftarrow} p(X/w_2) P(w_2) \stackrel{\rightarrow}{\rightarrow} X \in \begin{cases} w_1 \\ w_2 \end{cases}$$
 (3.5)

$$\mathcal{L}(\mathbf{X}) = \frac{\mathbf{p}(\mathbf{X}/\mathbf{w}_1)}{\mathbf{p}(\mathbf{X}/\mathbf{w}_2)} \geq \frac{\mathbf{P}(\mathbf{w}_2)}{\mathbf{P}(\mathbf{w}_1)} \rightarrow \mathbf{X} \in \begin{cases} \mathbf{w}_1 \\ \mathbf{w}_2 \end{cases}$$
(3.6)

h(X) can then be written as:

$$h(X) = -in (f(X)) = \frac{1}{2} (X - M_1)^T f_1^{-1} (Y - M_1) - \frac{1}{2} (X - M_2)^T f_2^{-1} (X - M_2)$$

$$+ \frac{1}{2} in \left| \frac{f_1}{f_2} \right| \geq \frac{1n}{P(w_1)} + K f_2^{-1} \left(\frac{w_2}{w_1} - \frac{w_2}{w_1} \right)$$

In practice, since M_i and Σ_i are replaced by \widetilde{M}_i and $\widetilde{\Sigma}_i$, h(X) becomes (after moving $\ln P(w_1)/P(w_2)$ to the L.H.S.): $\hat{h}(X) = \frac{1}{2}(X - \widetilde{M}_1)^T \hat{\Sigma}_1^{-1} (X - \widetilde{M}_1) - \frac{1}{2}(X - \widetilde{M}_2)^T \hat{\Sigma}_2^{-1} (X - \widetilde{M}_2)$

$$\frac{1}{2} \ln \left| \frac{\hat{r}_{1}}{\hat{r}_{2}} - \ln \frac{\hat{P}(w_{1})}{\hat{P}(w_{2})} \right| \geq 0 \rightarrow X \epsilon \begin{cases} w_{2} & (3.8) \\ w_{1} & (3.8) \end{cases}$$

The Bayes test for minimum error reduces then to looking at the value of $\hat{h}(X)$, assigning measurements with positive values to class 2, and measurements with negative values to class 1.

Note that h(X) is a one-dimensional random variable. The problem then is to know, or estimate, the probability density function of $\hat{h}(X)$. Once that is known, the probability of error can be obtained by carrying out a scalar integration. Figure 3.1 shows the probability density functions for h(X) given either class 1 or 2.

The probability of error can be calculated as:

 $\varepsilon = p(error) = p(error/w_1)P(w_1) + p(error/w_2)P(w_2)$ (3.9)

Discriminant function:



Figure 3.1 Probability Density Functions of h(X/w_i) and The Probability of Error.

Let the domain or decision space of X be divided into regions Γ_1 and Γ_2 . Then, if a sample belongs to w_1 , an error occurs whenever $X \in \Gamma_2$. Similarly, if a sample belongs to w_2 , an error occurs whenever $X \in \Gamma_1$. Thus,

$$\varepsilon = \Gamma \left(X \varepsilon \Gamma_2 / w_1 \right) P(w_1) + P(X \varepsilon \Gamma_1 / w_2) P(w_2)$$
(3.10)

In terms of the probability density functions of $\hat{h}(X/w_{_{f}})$, this becomes:

$$\varepsilon = P(w_1) \int_0^\infty p(h/w_1) dh + P(w_2) \int_0^\infty p(h/w_2) dh$$

$$= e_1 + e_2$$
(3.11)

The probability of error is then the area under the two curves in Figure 3.1 multiplied by the prior probabilities. The objective is to develop an algorithm which will approximate the class-conditional probability of $\hat{h}(X)$, and hence, the probability of error.

3.2 Performance Estimator

Fukunaga and Krile (64) developed an algorithm that approximates $\hat{h}(X)$. This algorithm assumes there are twoclass multivariate normal distributions, and was tested using one eight-dimensional simulated data set.

The algorithm, however, assumes the training samples are enough to reasonably estimate the true statistics of the distributions, and hence does not take into account the Hughes phenomenon. Put in other words, in situations where the training samples are few and do not reflect the true statistics of the distributions, the algorithm will treat the statistics obtained from the training samples as a "perfect" estimation of some "wrong" distributions, when in fact they are an "imperfect" estimation of the true statistics.

It is this algorithm, proposed by Fukunaga and Krile, that we will use and modify to take into account the Hughes phenomenon. Therefore, it seems appropriate to explain the algorithm in detail, and then discuss the modifications made to it.

3.2.1 The Normal Assumption

Looking at equation (3.8), since h(X) is a quadratic function in general of a normal random variable X, it cannot itself in general be normally distributed. However, in the case where $\Sigma_1 = \Sigma_2$, $\hat{h}(X)$ becomes a linear function of X and hence is normally distributed.

In most cases, however, $\Sigma_1 \neq \Sigma_2$. Fukunaga and Krile still tried to assume that $\hat{h}(X)$ is normally distributed.

An algorithm was developed and tested in this research under the assumption that $\hat{h}(X)$ is normally distributed (although $\Sigma_1 \neq \Sigma_2$) but results showed it to be a very poor approximation of the probability of error and hence it was not further analyzed.

3.2.2 The Modified Gamma Distribution Assumption: Fukunaga and Krile Version

Consider $\hat{h}(X)$ as given by equation (3.8). Applying the simultaneous diagonalization technique described earlier, $\hat{\Sigma}_1$ is transformed to the identity matrix I, and $\hat{\Sigma}_2$ is transformed to a diagonal matrix Λ . The transformation matrix is denoted A^T , or the transpose of the eigenvector matrix A.

Without losing generality, we assign the origin of the coordinate system such that:

$$\hat{m}_1 = 0$$
 and $\hat{m}_2 = \hat{M}_1 - \hat{M}_2$ (3.12)

With X_{EW_1} , $\hat{h}(X)$ can be written as another function of Y, where $Y=A^T X$, as follows:

$$\hat{h}(Y/w_{1}) = Y^{T} Y - (Y-\hat{D})^{T} \hat{\Lambda}^{-1} (Y-\hat{D}) + \ln \left| \frac{\hat{\Sigma}_{1}}{\hat{\Sigma}_{2}} \right| \\ - 2 \ln \frac{\hat{P}(w_{1})}{\hat{P}(w_{2})}$$
(3.13)
where $\hat{D} = A^{T}\hat{m}_{2}$.

Since the features are now uncoupled, this can be written as:

$$\hat{\mathbf{h}}(\mathbf{Y}/\mathbf{w}_{1}) = \sum_{i=1}^{p} (\mathbf{y}_{i}^{2} - \frac{1}{\hat{\lambda}_{i}} (\mathbf{y}_{i} - \hat{\mathbf{d}}_{i})^{2} - \ln \hat{\lambda}_{i}) - 2 \ln \frac{\hat{\mathbf{P}}(\mathbf{w}_{1})}{\hat{\mathbf{P}}(\mathbf{w}_{2})}$$

$$= \sum_{i=1}^{p} ((1 - \frac{1}{\hat{\lambda}_{i}}) (\mathbf{y}_{i} + \frac{\hat{\mathbf{d}}_{i}}{\hat{\lambda}_{i} - 1})^{2} - (\frac{\hat{\mathbf{d}}_{i}^{2}}{\hat{\lambda}_{i} - 1} + \ln \hat{\lambda}_{i}))$$

$$- 2 \ln \frac{\hat{\mathbf{P}}(\mathbf{w}_{1})}{\hat{\mathbf{P}}(\mathbf{w}_{2})}$$
(3.14)

where p is the number of dimensions. \hat{d}_{i} is the ith element of vector D.

Now, we have $\hat{h}(Y/w_1)$ in terms of p independent Gaussian random variables y_1 , each of which has zero mean and unit variance with respect to class w_1 .

Defining a new transformed variable Z and a transformed difference- of-means vector $\hat{\nu}$ as follows:

$$Z = (\hat{\Lambda}^{-\frac{1}{2}} A^{T}) (X - \hat{m}_{2})$$
(3.15)

$$\hat{v} = (\hat{\Lambda}^{-\frac{1}{2}} A^{T}) \hat{m}_{2} = \hat{\Lambda}^{-\frac{1}{2}} \hat{D}$$
 (3.16)

 $\hat{h}(X/w_2)$ can be expressed as a function of the new variable Z and \hat{v} by substituting (3.15) and (3.16) into (3.8) as follows:

$$\hat{\mathbf{n}}(z/w_2) = (z+\hat{\mathbf{v}})^T \hat{\mathbf{n}}(z+\hat{\mathbf{v}}) - z^T z + \ln \left| \frac{\hat{\mathbf{z}}_1}{\hat{\mathbf{z}}_2} \right| - 2 \ln \frac{\hat{\mathbf{p}}(w_1)}{\hat{\mathbf{p}}(w_2)}$$
(3.17)

Again, since the features are uncoupled, we can write $\hat{h}(Z/w_2) = \sum_{i=1}^{p} (\hat{\lambda}_i (z_i + \hat{\nu}_i)^2 - z_i^2 - \ln \hat{\lambda}_i) - 2 \ln \frac{\hat{P}(w_1)}{\hat{P}(w_2)}$ $= \sum_{i=1}^{p} ((\hat{\lambda}_i - 1) (z_i + \frac{\hat{\lambda}_i^{k_2} \hat{d}_i}{\hat{\lambda}_i - 1})^2 - \frac{\hat{d}_i}{\hat{\lambda}_i - 1} + \ln \hat{\lambda}_i))$ $- 2 \ln \frac{\hat{P}(w_1)}{\hat{P}(w_2)} \qquad (3.18)$

Again, we have an expression in terms of p independent Gaussian variables z_i , each of which has zero mean and unit variance.

Next, we define the following quantities for convenience:

- $a_{1i} = 1 1/\hat{\lambda}_i$ (3.19)
- $b_{1i} = \hat{d}_i / (\hat{\lambda}_i 1)$ (3.20)
- $a_{2i} = \hat{\lambda}_i 1$ (3.21)

$$b_{2i} = \hat{\lambda}_{i}^{k_{2}} \hat{d}_{i}^{/(\hat{\lambda}_{i}-1)}$$
(3.22)

$$C = \sum_{i=1}^{P} (\ln \hat{\lambda}_{i} + \hat{d}_{i}/(\hat{\lambda}_{i}-1) + 2 \ln \hat{P}(w_{1})/\hat{P}(w_{2})$$
(3.23)

Substituting equations (3.19)-(3.23) back into equations (3.14) and (3.18), we get:

$$\hat{h}(Y/w_1) = \sum_{i=1}^{p} (a_{1i} (y_i + b_{1i})^2) - C$$
 (3.24)

$$\hat{h}(Z/w_2) = \sum_{i=1}^{p} (a_{2i} (z_i + b_{2i})^2) - C$$
 (3.25)

Referring from now on to Y and Z as ξ , and to y_i and z_i as ξ_i , we find that $\hat{h}(\hat{\xi}/w_1)$ and $h(\hat{\xi}/w_2)$ have the same functional form, except for the values of a_{1i} , b_{1i} , a_{2i} , and b_{2i} .

Theorem 3.1

If $X = (x_1, \dots, x_p)$ where the x_i are a sample from a Normal(0, σ^2) population, then the random variable V = $\sum_{\substack{p \\ i=1}}^{p} \chi_p^2 / \sigma^2$ has a χ_p^2 , or chi-square, distribution. <u>Proof</u>:

See (67), p. 16.

Theorem 3.2

If s_1, \ldots, s_p are independent random variables, then the density of their sum $s_1 + s_2 + \ldots + s_p$ equals the convolution of their respective densities. <u>Proof</u>

```
See (68), p. 189.
```

Examining equations (3.24) and (3.25), shows that the density functions of $\hat{h}(\xi/w_1)$ and $\hat{h}(\xi/w_2)$ can be obtained by convolving the densities of p non-central (because of the b_{1i} and the b_{2i} terms) χ^2 variables having multiplicative constants a_{1i} and a_{2i} , and adding a shift parameter C.

The density of $\hat{h}(\xi)$ is divided into three parts:

$$V_{kr} = \sum_{\substack{ki \geq 0 \\ a_{ki} \geq 0}} a_{ki} (\xi_{ki} + b_{ki})^2 \text{ for } a_{ki} \geq 0 \quad (3.26)$$

$$V_{ks} = \Sigma = a_{kj} (\xi_{kj} + b_{kj})^2 \text{ for } a_{kj} < 0 \qquad (3.27)$$

$$C = \sum_{i=1}^{p} (\ln \hat{\lambda}_{i} + \hat{d}_{i}/(\hat{\lambda}_{i}-1) + 2 \ln \hat{P}(w_{1})/\hat{P}(w_{2})$$
(3.28)

$$(p = p_{kr} + p_{ks})$$
 (k = 1,2)

The density function of V_{kr} , $p_{kr}(h)$, is the convolution of p_{kr} densities of squared Gaussian variables having multiplicative constants. All p_{kr} densities lie above the positive h axis with $a_{ki} \ge 0$. Similarly, the density function of V_{ks} , $p_{ks}(h)$, is the convolution of p_{ks} densities of squared Gaussian variables with multiplicative constants. All p_{ks} densities lie on the negative h axis with $a_{ki} \le 0$.

A gamma density function is given by:

$$g_{p+\lambda} = \lambda^{p} x^{p-1} e^{-\lambda x} / \Gamma(p)$$
 (3.29)

Let k be a positive integer. With p=1/2k, and $\lambda = 1/2$, the gamma density $g(p,\lambda)$ is referred to as the chi-squared density with k degrees of freedom. (See (67), p.13).

Theorem 3.3

If X_1, \ldots, X_n are independent random variables with gamma distributions $(p_1, \lambda), \ldots, (p_n, \lambda)$, then $Y = X_1 + \ldots + X_n$ has a gamma distribution $(p_1 + \ldots + p_n, \lambda)$. Proof

See (67). p. 15.

Since what we have is the summation of chi-squared random variables (special form of a gamma distribution), both $p_{kr}(h)$ and $p_{ks}(-h)$ ($p_{ks}(h)$ reflected to the positive side) can be reasonably approximated by a general gamma form, especially for large n_{kr} and n_{ks} , as follows:

$$g(h) = \begin{cases} \frac{h^{\alpha} e^{-h/\ell}}{\beta^{\alpha+1} r(\alpha+1)} & h \ge 0\\ 0 & h \le 0 \end{cases}$$
(3.30)

The parameters α and β can be determined so that the mean n and the variance σ^2 of the "true" distribution match those of the approximation.

Next, we calculate the expected values n_{kr} and n_{ks} of V_{kr} , and V_{ks} , and the variances σ_{kr}^2 and σ_{ks}^2 . $v_{kr} = \sum_{\substack{ki \geq 0}}^{p_{kr}} a_{ki} (\xi_{ki} + b_{ki})^2 \qquad a_{ki \geq 0}^2$ $= \sum_{\substack{k_{1} \\ a_{k_{1}} \geq 0}}^{p_{k_{1}}} a_{k_{1}} (\xi_{k_{1}}^{2} + 2b_{k_{1}}\xi_{k_{1}} + b_{k_{1}}^{2})$ $E(V_{kr}) = \hat{n}_{kr} = \sum_{a_{1,4} \ge 0}^{p_{kr}} (1 + 0 + b_{ki}^2)$ or, $\hat{p}_{kr} = \sum_{ki}^{p} a_{ki} (1 + b_{ki}^2)$ for $p_{kr}(h)$ $a_{ki} \ge 0$ (3.31)

 $(\xi_{ki}$ has zero mean and unit variance) Similarly, P.

$$\hat{n}_{ks} = \hat{\Sigma} = a_{kj} (1 + b_{kj}^2) \text{ for } p_{ks}(h)$$
 (3.32)
 $a_{kj} = 0$

$$E(V_{kr}^{2}) = E(\sum_{a_{ki}, a_{kj}}^{p_{kr}} \sum_{a_{kj}, a_{kj}}^{p_{kr}} \sum_{a_{ki}, a_{kj}}^{p_{ki}} \sum_{a_{ki}, a_{kj}}^{p_{ki}} \sum_{a_{ki}, a_{kj}}^{p_{ki}} \sum_{a_{ki}, a_{kj}}^{p_{ki}} \sum_{a_{ki}, a_{ki}}^{p_{ki}} \sum_{a_{ki},$$

(The zero term comes because ξ_{ki} is independent from ξ_{kj} and hence they are mutually orthogonal as $E(\xi_{ki}) = 0$)

$$= \sum_{\substack{ki=0}}^{p_{kr}} a_{ki}^{2} (3 + 6 b_{ki}^{2} + b_{ki}^{4})$$
(3.33)

where
$$E(\xi_{ki}^{n}) = \begin{cases} 1.3. \dots (n-1) & \text{for } n \text{ even} \\ 0 & \text{for } n \text{ odd} \end{cases}$$

$$E^{2} (V_{kr}) = \sum_{\substack{a_{ki} \geq 0 \\ a_{ki} \geq 0}}^{p_{kr}} a_{ki}^{2} (1 + b_{ki}^{2}) + 0$$
$$= \sum_{\substack{a_{ki} \geq 0 \\ a_{ki} \geq 0}}^{p_{kr}} a_{ki}^{2} (1 + 2 b_{ki}^{2} + b_{ki}^{4})$$
(3.34)

$$v_{ar} (v_{kr}) = \hat{\sigma}_{kr} = E(v_{kr}^2) - E^2 (v_{kr})$$

= 2 $\sum_{ki=0}^{p_{kr}} a_{ki}^2 (1 + 2b_{ki}^2) \text{ for } p_{kr}(h)$ (3.35)
 $a_{ki}^2 = 0$

Similarly,

$$\hat{\sigma}_{k5} = 2 \sum_{\substack{k = 0 \\ k \neq 0}}^{p_{k5}} a_{kj}^{2}(1 + 2 b_{kj}^{2}) \text{ for } p_{ks}(h) \quad (3.36)$$

For a random variable h, which has a gamma distribution with parameters α and β , (See equation (3.30)), then

$$E(h) = (\alpha + 1)\beta$$
 Var (h) = $(\alpha + 1)\beta^2$ (3.37)
(See (67), p. 44)

Therefore, a_{kr} , a_{ks} , β_{kr} , β_{ks} , can be calculated as:

$$\alpha_{kr} = (\hat{n}_{kr}^2 / \hat{\sigma}_{kr}^2) - 1$$
 (3.38)

$$\alpha_{ks} = (\hat{\eta}_{ks}^2 / \hat{\sigma}_{ks}^2) - 1$$
 (3.39)

$$\beta_{\mathbf{k}\mathbf{r}} = \hat{\sigma}_{\mathbf{k}\mathbf{r}}^2 / \hat{\eta}_{\mathbf{k}\mathbf{r}}$$
(3.40)

$$\beta_{ks} = \sigma_{ks}^2 / \eta_{ks}$$
(3.41)

The density function $p(h/w_i)$, i=1,2, which is our final goal, is then the convolution of two gamma densities with a constant shift: one is distributed on the positive side of the h-axis, and the other on the negative side.

However, the convolution of these two gamma densitities is hard to obtain in an explicit mathematical expression, because in general, α is not an integer. Since we do not favor a numerical integration technique for calculating the error rate, a "modified " gamma distribution is proposed as follows:

$$g'(h) = \begin{cases} \frac{(h-c)^{\gamma} e^{-(h-c)/\delta}}{\delta^{\gamma+1} \Gamma(\gamma+1)} & \text{for } h \ge c \\ 0 & \text{for } h < c \end{cases}$$
(3.42)

 $\gamma = 0 \text{ or } 1$

In other words, Gamma density curves are roughly categorized into two types: one is $exp(-h/\beta)$, and the other is h $exp(-h/\beta)$, depending on whether α obtained by (3.38) or (3.39) is larger than or smaller than a threshold value of 0.35. (The threshold value of 0.35 is a compromise value, chosen in an attempt to match the maximum value and location of the maximum value of the gamma density to the modified gamma approximation. It is further explained in (64)).

The procedure proposed by Fukunaga and Krile, then, is as follows:

1) Calculate $\hat{\eta}_{kr}$, $\hat{\eta}_{ks}$, $\hat{\sigma}_{kr}^2$, $\hat{\sigma}_{ks}^2$ from equations (3.31),(3.32),(3.35), and (3.36)

- 2) Calculate a_{kr} and a_{ks} form equations (3.38) and (3.39).
- 3) $\gamma_{kr} = 0$ if $\alpha_{kr} < 0.35$, and $\gamma_{kr} = 1$ if $\alpha_{kr} > 0.35$. Similarly for γ_{ks} .
- 4) Calculate δ_{kr} , $\delta_{ks'}$, and c_{kr} , c_{ks} by the following equations: (modified forms of equations (3.38)-(3.41))

$$\gamma_{kr} = \frac{(\hat{n}_{kr} - c_{kr})^2}{\hat{\sigma}_{kr}^2} - 1 \qquad (3.43)$$

$$\gamma_{ks} = \frac{(\hat{\eta}_{ks} - c_{ks})^2}{\hat{\sigma}_{ks}^2} - 1$$
 (3.44)

$$\delta_{kr} = \hat{\sigma}_{kr}^2 / (\hat{n}_{kr} - c_{kr})$$
 (3.45)

$$\delta_{ks} = \hat{\sigma}_{ks}^2 / (\hat{n}_{ks} - c_{ks})$$
 (3.46)

Equations (3.43)-(3.46) are the same as (3.38)-(3.41), except for the shift of the mean c_{kr} or c_{ks} .

The convolution of $p_{kr}(h)$ and $p_{ks}(h)$, $p_{k}^{*}(h)$, k=1,2, can be obtained as an explicit expression. The result is : (See (64) for details)

$$P_{k}^{*}(t) = \begin{cases} \frac{\delta_{ks}^{\gamma_{kr}}}{(\delta_{kr} + \delta_{ks})^{\gamma_{kr+1}}} \begin{bmatrix} \frac{t}{\delta_{ks}} + \frac{(\gamma_{kr} + \gamma_{ks})^{\delta_{kr}}}{\delta_{kr} + \delta_{ks}} \end{bmatrix}^{\gamma_{ks}} e^{t/\delta_{ks}} \\ for t \leq 0 \\ \frac{\delta_{kr}^{\gamma_{ks}}}{(\delta_{kr} + \delta_{ks})^{\gamma_{ks+1}}} \begin{bmatrix} \frac{t}{\delta_{kr}} + \frac{(\gamma_{kr} + \gamma_{ks})^{\delta_{ks}}}{\delta_{kr} + \delta_{ks}} \end{bmatrix}^{\gamma_{kr}} e^{-t/\delta_{kr}} \\ for t \geq 0 \\ (3.47) \end{cases}$$

Defining the distance d as

$$d_k = C - (c_{kr} - c_{ks})$$
 (3.48)

We can find e_1 by integrating $p_1^*(t)$ form d_1 to ∞ , and e_2 by integrating $p_2^*(t)$ from $-\infty$ to d_2 . The term d_k brings the shift parameter C back ino the picture, and also accounts for the displacement of the (h/w_k) approximations by c_{kr} and c_{ks} . In general,

$$D^{\star}(d_{k}) = \int_{-\infty}^{d_{k}} p_{k}^{\star}(t) dt =$$

$$\begin{pmatrix} \frac{\delta_{ks}}{\delta_{kr} + \delta_{ks}} \end{pmatrix}^{\gamma_{kr+1}} \begin{bmatrix} \frac{-d_{k}}{\delta_{ks}} + 1 + \frac{(\gamma_{kr} + \gamma_{ks})\delta_{kr}}{\delta_{kr} + \delta_{ks}} \end{bmatrix}^{\gamma_{ks}} e^{d_{k}/\delta_{ks}}, \quad d_{k} \leq 0$$

$$1 - \left(\frac{\delta_{kr}}{\delta_{kr} + \delta_{ks}} \right)^{\gamma_{ks+1}} \begin{bmatrix} \frac{d_{k}}{\delta_{kr}} + 1 + \frac{(\gamma_{kr} + \gamma_{ks})\delta_{ks}}{\delta_{kr} + \delta_{ks}} \end{bmatrix}^{\gamma_{kr}} e^{-d_{k}/\delta_{kr}}, \quad d_{k} \geq 0$$

(3.49)

where $\Gamma^{\dagger}(d_k)$ is the approximation for $Prob(h/w_k \leq 0)$. Thus, the approximated values of recognition errors are:

$$e_1 = \hat{P}(w_1) (1 - D^{\pi}(d_1))$$
 (2.50)

$$e_2 = \hat{P}(w_2) (D^*(d_2))$$
 (3.51)

3.2.3 Proposed, Modified Algorithm

Figure 3.2 shows a flowchart of Fukunaga's and Krile's algorithm. The algorithm assumes that the training statistics are an accurate representation of the true statistics of the two distributions. This being the case, the probability of correct classification that the algorithm projects is monotonically non-decreasing as a function of dimensionality. It is this drawback in the algorithm that we are trying to correct such that the algorithm would take into account the number of samples used for training.

Looking back at the calculation of the parameters of the modified gamma distribution, we see that all of them depend on two parameters, \hat{n}_k and $\hat{\sigma}_k$, or the mean and variance of h. If these parameters are inaccurate, then all of the other parameters will be affected.



Figure 3.2 A Flowchart of Fukunaga and Krile's Algorithm.

*

We propose to look it the way these parameters, particularly $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$, are distributed as a function of the number of training samples. We then want to incorporate that information in our estimation of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$, such that the algorithm has a more realistic picture of what the training samples represent.

Estimating the probability density function of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ is by no means an easy task. For the amount of information that we have, such an estimation is very involved and impractical. A discussion of the difficulties one faces in attempting such an estimation is found in Appendix B.

We propose instead to look at the variances of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$, and then incorporate that information in our estimation of these parameters.

Let us look at $\hat{\sigma}_1^2$, (Var (\hat{h}/w_1)) and $\hat{\sigma}_2^2$, (Var (\hat{h}/w_2)). From equation (3.35), (or (3.36)):

$$\hat{\sigma}_{1}^{2} = 2 \sum_{i=1}^{p} a_{1i}^{2} (1 + 2 b_{1i}^{2})$$
(3.35)

Substituting for a_{1i} and b_{1i} by their values from (3.19) and (3.20) in (3.35), we get:

$$\hat{\sigma}_{1}^{2} = 2 \sum_{i=1}^{p} (1 - 1/\hat{\lambda}_{i})^{2} (1 + 2 \hat{d}_{i}^{2}/(\hat{\lambda}_{i}^{-1})^{2})$$
(3.52)

After multiplying, this reduces to:

$$\hat{\sigma}_{1}^{2} = 2 \frac{P}{1-1} (1 - 2/\hat{\lambda}_{1} + (2 \hat{a}_{1}^{2} + 1) / \hat{\lambda}_{1}^{2}) \qquad (3.53)$$
In matrix form, this can be written as:

$$\hat{\sigma}_{1}^{2} = 2 (tr (I - \hat{\lambda}^{-1})^{2} + 2 \hat{D}^{T} (\hat{\lambda}^{-1})^{2} \hat{D}) \qquad (3.54)$$
Or in terms of the original distributions:

$$\hat{\sigma}_{1}^{2} = 2 (tr (I - \hat{\Sigma}_{2}^{-1} \hat{\Sigma}_{1})^{2} + 2 \hat{m}_{2}^{T} \hat{\Sigma}_{2}^{-1} \hat{\Sigma}_{1} \hat{\Sigma}_{2}^{-1} \hat{m}_{2}) \qquad (3.55)$$
(See (64)).
Similarly,

$$\hat{\sigma}_{2}^{2} = 2 \frac{P}{1-1} a_{21}^{2} (1 + 2 b_{21}^{2})$$

$$= 2 \frac{P}{1-1} (\hat{\lambda}_{1} - 1)^{2} (1 + 2 \hat{\lambda}_{1} \hat{d}_{1}^{2} / (\hat{\lambda}_{1} - 1)^{2})$$

$$= 2 \frac{P}{1-1} (\hat{\lambda}_{1}^{2} + 2 (\hat{d}_{1} - 1) \hat{\lambda}_{1} + 1) \qquad (3.56)$$
In matrix form, $\hat{\sigma}_{2}^{2}$ can be written as:

 $\hat{\sigma}_{2}^{2} = 2 (tr (\hat{\Lambda} - I)^{2} + 2 \hat{D}^{T} \hat{\Lambda} \hat{D})$ (3.57)

Or, in terms of the original distributions:

$$\hat{\sigma}_2^2 = 2 (tr (\hat{\Sigma}_1^{-1} \hat{\Sigma}_2 - I)^2 + 2 \hat{m}_2^T \hat{\Sigma}_1^{-1} \hat{\Sigma}_2 \hat{\Sigma}_1^{-1} \hat{m}_2)$$
 (3.58)

(See (64)).

In order to calculate the variances of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$, we make the following assumptions:

- 1. The original and transformed means, \hat{M}_1, \hat{M}_2 , and \hat{D} are assumed to be constant. Experience has shown that one can approximate first- order statistics with a relatively few number of training samples.
- 2. $\hat{\Sigma}_1$ and $\hat{\Sigma}_2$ are independent. This is to say that we will ignore any relationships that might exist between the two classes.

Having assumed the above, the results are: (See Appendix C for the complete derivation)

$$\operatorname{Var}(\hat{\mathfrak{o}}_{1}^{2}) = 4 \frac{p}{\sum_{i=1}^{p}} \left(\frac{2}{\lambda_{1}^{2}} \left(\frac{4}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{1}n_{2}} \right) - \frac{4}{\lambda_{1}^{3}} \left(\frac{4}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{2}} + \frac{8}{n_{2}} + \frac{4}{\lambda_{1}^{2}} + \frac{4}{n_{2}} + \frac{4}{n_{2}} + \frac{8}{n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}} + \frac{4}{n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}} + \frac{4}{n_{2}^{2}} + \frac{4}{n_{1}^{2}n_{2}^{2}} + \frac{4}{n_{1}^{2}$$

$$+ \frac{1920}{n_1^2 n_2^2} + \frac{576}{n_1^3 n_2} + \frac{576}{n_2^3 n_1} + \frac{2112}{n_1^2 n_2^2} + \frac{2112}{n_1^3 n_2^2} + \frac{2304}{n_1^3 n_2^3} + 4d_1^2 \left(\frac{4}{n_1} + \frac{8}{n_2} + \frac{8}{n_1^2 n_2^2} + \frac{40}{n_1^2 n_2^2} + \frac{256}{n_1^2 n_2^2} + \frac{96}{n_1^2 n_2^2} + \frac{48}{n_2^3} + \frac{288}{n_1 n_2^3} + \frac{352}{n_1^2 n_2^2} + \frac{384}{n_1^2 n_2^3}\right)$$

$$+ 4d_1^4 \left(\frac{2}{n_1} + \frac{8}{n_2} + \frac{40}{n_2^2} + \frac{24}{n_1 n_2} + \frac{48}{n_2^3} + \frac{88}{n_2^3} + \frac{96}{n_1 n_2^3}\right)$$

$$(3.59)$$

$$Var(\hat{\sigma}_2^2) = 4 \frac{p}{1-1} \left[\lambda_1^4 \left(\frac{8}{n_1} + \frac{8}{n_2} + \frac{128}{n_1 n_2^3} + \frac{40}{n_1^2 n_2^2} + \frac{49}{n_1^2 n_2^2} + \frac{48}{n_1^3} + \frac{48}{n_2^3} + \frac{512}{n_1^2 n_2^2} + \frac{576}{n_1^2 n_2^2} + \frac{576}{n_1^2 n_2^2} + \frac{2112}{n_1^2 n_2^2} + \frac{40}{n_1^2 n_2^3} \right) + 4\lambda_1^3 \left(d_1^2 \left(\frac{8}{n_1} + \frac{4}{n_2} + \frac{8}{n_1^2} + \frac{40}{n_1^2 n_2^2} + \frac{2112}{n_1^2 n_2^2} + \frac{2304}{n_1^2 n_2^3} \right) + 4\lambda_1^3 \left(d_1^2 \left(\frac{8}{n_1} + \frac{4}{n_2} + \frac{8}{n_2^2} + \frac{40}{n_1^2 n_2^2} + \frac{256}{n_1^2 n_2} + \frac{96}{n_2^2 n_1} + \frac{48}{n_1^3} + \frac{288}{n_2^2 n_1} + \frac{352}{n_1^2 n_2^2} + \frac{384}{n_2^2 n_1^2} \right) \right)$$

$$- \left(\frac{4}{n_1} + \frac{4}{n_2} + \frac{8}{n_1^2} + \frac{8}{n_2^2} + \frac{32}{n_1^2 n_2} + \frac{32}{n_1^2 n_2} + \frac{48}{n_1^2 n_2^2} + \frac{48}{n_1^2 n_2} + \frac{64}{n_1^2 n_2^2} \right) \right) + 2\lambda_1^2 \left(\left(\frac{4}{n_1} + \frac{4}{n_2} + \frac{8}{n_1^2} + \frac{8}{n_2^2} + \frac{32}{n_1^2 n_2} + \frac{48}{n_1^2 n_2^2} + \frac{48}{n_1^2 n_2} + \frac{48}{n_1^2 n_2^2} + \frac{64}{n_1^2 n_2^2} \right) \right) \right)$$

$$- 4d_1^2 \left(\frac{2}{n_2} + \frac{4}{n_1} + \frac{12}{n_1^2 n_2} + \frac{8}{n_1^2} + \frac{16}{n_1^2 n_2} \right) \right) \right]$$

$$(3.60)$$

Note that $Var(\hat{\sigma}_1^2)$ and $Var(\hat{\sigma}_2^2)$ are inversely proportional to the number of training samples used to estimate the statistics of classes 1 and 2, and directly proportional to the number of dimensions. In other words, as the number of training samples increases, the variances of our

6 ï
estimates of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ decrease, as expected. Also, as the number of dimensions is increased, the variances of the estimates increase.

Since we do not have the probability density functions of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$, we want to think of a reasonable way to incorporate the effect of the number of training samples into our estimation of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$. We claim that a better estimation of the true variances σ_1^2 and σ_2^2 consists of our estimation of these variances, $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$, plus some multiplicative factor of the standard deviations of these estimates, namely the square roots of $\operatorname{Var}(\hat{\sigma}_1^2)$ and $\operatorname{Var}(\hat{\sigma}_2^2)$, that were calculated above.

This multiplicative factor was chosen empirically. Experimental results in Chapter 4 show that the variance of the probability of error generally increases with increasing dimensionality, especially in the presence of a very limited training data set. Results also show that the probability of error is inversely proportional to the number of training samples. Moreover, it is very sensitive to the number of training samples in the cases where that number is not much greater than the number of dimensions.

Based on the above observations, the following empirical formula for the multiplicative factor was used:

Part days in a will be a

M.F. =
$$2 p^2 / (n_1 \cdot n_2)$$
 (3.61)

where p is the number of dimensions

n₁ and n₂ are as before.

The new procedure to calculate law probability of error, becomes as follows:

- 1) Calculate \hat{n}_{kr} , \hat{n}_{ks} , $\hat{\sigma}_{kr}$, $\hat{\sigma}_{ks}$, from equations (3.31), (3.32), (3.35), and (3.36)
- 2) Update $\hat{\sigma}_{kr}^2$ and $\hat{\sigma}_{ks}^2$ as follows: $\hat{\sigma}_{kr}^2$ (new) = $\hat{\sigma}_{kr}^2$ (old) + $(2p^2/n_1 \cdot n_2) \cdot (Var(\hat{\sigma}_{kr}^2))^{\frac{1}{2}}$ $\hat{\sigma}_{ks}^2$ (new) = $\hat{\sigma}_{ks}^2$ (old) + $(2p^2/n_1 \cdot n_2) \cdot (Var(\hat{\sigma}_{ks}^2))^{\frac{1}{2}}$
- 3) $\gamma_{kr} = 1$ if $\alpha_{kr} \ge 0.35$, and $\gamma_{kr} = 0$ if $\alpha_{kr} < 0.35$. Similarly for γ_{ks} .
- 4) Calculate δ_{kr} , δ_{ks} , and c_{kr} , c_{ks} , from equations ((3.43) - (3.46)).
- 5) Calculate p_k^* (t) and $D^*(d_k)$ from equations (3.47) and (3.49).
- 6) Calculate the probability of error from equations (3.59) and (3.60).

We are ready now to proceed to Chapter 4, where several experimental results are shown.

CHAPTER 4

EXPERIMENTAL RESULTS

4.1 Introduction

Some results on feature selection techniques will be presented first. Next, several experimental results illustrating the Hughes phenomenon are shown. Results comparing probabilities of error predicted by the proposed algorithm as a function of dimensionality as compared to experimental observations are then presented for aircraft and Landsat data. Results are obtained for both real and simulated data. Finally, two binary tree classification procedures that make use of the algorithm are presented to illustrate the usefulness of the procedure.

The Bayesian decision rule with assumptions of 0-1 loss function, equal a priori probabilities, and multivariate normal distributions is used as the decision rule in all experiments when classification is involved.

Detailed training and test field descriptions for all the experiments conducted are found in Appendix F.

4.2 Experiments on Feature Selection Techniques

In this section, some experiments on different feature selection techniques are presented. The purpose of conducting these experiments is to choose an effective feature selection technique, particularly when dealing with a small number of training samples.

Experiment 4.1

Two classes of wheat and corn are selected from multispectral scanner (hereafter referred to as MSS or aircraft) data of the 1971 Flightline 210 from the Corn Blight Watch Experiment, and classified. The data was collected on August 13, 1971. Part of the selected data is used for training and a much larger portion is used for testing. The number of features used for classification varies from one to twelve, and the number of training samples for each class is chosen such that it is much higher than the number of features (265 samples for wheat, 569 samples for corn). A principle components (Karhunen-Loeve) transformation is applied to the data, and then three feature selection techniques are compared:

 In the first feature selection method, the features are ordered according to the largest eigenvalues resulting from the K-L expansion. This method, referred to hereafter as the K-L ordering method,

assumes that the best feature is that which corresponds to the largest eigenvalue of the mixture covariance matrix of the whole data set, the second best corresponds to the second largest eigenvalue, ...etc. This ordering then imposes the condition that a feature subset with lower dimensionality is always a subset of another with higher dimensionality. The method then depends on the eigenvalues of the mixture covariance matrix, and ignores any among-class variabilities.

2) The second feature selection technique method is referred to as the Transformed Divergence method (13). The transformed divergence, D_T, is defined as follows:

$$D_{\rm T} = 2000 \, (1 - e^{-D/8})$$
 (4.1)

where D is the divergence of two normal distributions, and is defined as follows (12):

$$D = \frac{1}{2} \operatorname{tr} (\hat{\Sigma}_{1} - \hat{\Sigma}_{2}) (\hat{\Sigma}_{2}^{-1} - \hat{\Sigma}_{1}^{-1}) + \frac{1}{2} (\hat{M}_{1} - \hat{M}_{2})^{T} (\hat{\Sigma}_{1}^{-1} + \hat{\Sigma}_{2}^{-1}) (\hat{M}_{1} - \hat{M}_{2})$$
(4.2)

For a given dimensionality, the method chooses the feature subset with that dimensionality which gives the largest value of D_{T^*} Unlike the K-L method, a feature subset of lower dimensionality is not neces-

-

sarily a subset of another with higher dimensionality. This method is applied to the data after it has been K-L transformed.

3) The third feature selection technique method used is the Bhattacharyya distance (16), defined by equation (2.9). In this method, a simultaneous diagonalization technique is applied to the covariance matrices of the two classes (after a K-L transformation of the data), and the best feature is then selected as that which corresponds to the largest value of B as defined by equation (2.10). The second largest is that which corresponds to the second largest B, and so on. As in the K-L method, a feature subset of lower dimensionality is always a subset of one with higher dimensionality. The transpose of the eigenvector matrix obtained is then multiplied by the observation vectors to transform the data, the mean vectors and the covariance matrices are transformed. and the data classified.

Results are shown in Figure 4.1, which plots the recognition accuracy (P_{cc} \$) as a function of dimensionality. It is seen that of the three methods, the transformed divergence one gives the poorest performance. The K-L method is better, but the best method is that obtained from the Bhattacharyya ordering, which saturates at a very low dimension-

, and the state

- 1**/A**



Figure 4.1 Classification Results of Data in Experiment 4.1 Using Three Feature Selection Techniques.

ality. Note that as dimensionality increases, the three curves start approaching each other, until they all coincide when all features are used (The probability of error is invariant under any linear transformation).

Experiment 4.2

In this experiment, 20 samples each of wheat and corn are chosen randomly from the training samples of experiment 4.1. The test samples are the same in both experiments. Again. the same three feature selection techniques elaborated upon above are used. Classification results are shown in Figure 4.2. Unlike the results in experiment 4.1, the Bhattacharyya ordering here gives the poorest results. Further, it does not exhibit a peaking effect, an effect that is expected when working with such a small number of training samples. The transformed divergence ordering does much better and does exhibit a peaking effect. However, it has a lot of fluctuations. The K-L ordering, on the other hand, while giving slightly poorer results than transformed divergence at low dimensionality, is better than the other two techniques at high dimensionality and has less fluctuations.



Figure 4.2 Classification Results of Data in Experiment 4.2 Using Three Feature Selection Techniques.

Experiment 4.3

Another two classes, corn and forest, are selected from the same data set described in experiment 4.1. Again, 20 samples per class are chosen randomly from a larger set of training samples, and the three feature selection techniques are compared. Results appear in Figure 4.3

Again, we notice that the Bhattacharyya ordering does poorer than the other two techniques, and does not exhibit a peaking effect. Transformed divergence gives better results, but again has a lot of fluctuations. The K-L ordering is superior to both, and has less fluctuations.

It should be noted again that the K-L ordering we used is based over the full data set. It is dependent on the mixture covariance matrix of the full data set, and thus ignores any between class variabilities resulting from differences between class covariance matrices. Because it is always dependent on the full data set, the number of training samples used to estimate the mixture covariance matrix is almost always large, and hence a good estimate is obtained.

The Bhattacharyya ordering used, on the other hand, although it takes into account between class variabilities, depends heavily on the number of training samples used to estimate the individual covariance matrices of the classes at hand. Thus, as the number of training samples decreases,



Figure 4.3 Classification Results of Data in Experiment 4.3 Using Three Feature Selection Techniques.

poorer estimates of the covariance matrices are obtained, leading to poorer transformations.

It appears that the transformation obtained from the simultaneous diagonalization technique is very sensitive to the number of training samples used to estimate the statistics of the classes at hand. While it produces superior results when there are enough samples, it fails to do so when the training samples are limited.

Indeed, Wu (50) published results in which he showed that the divergence criterion breaks down when the number of training samples is small, and no longer is an effective predictor of accuracy.

The K-L ordering, while ignoring the among-class variabilities in the scene, is only dependent on the number of data points in the data set used to approximate the mixture covariance matrix, but is otherwise independent of the number of training samples used. Thus, while sacrificing the information we get about the variability between classes in the set, experimental results show that this sacrifice is more than warranted when dealing with a small number of training samples. While not claiming that the K-L ordering gives the optimal results, we think it is a very effective procedure 10 the presence of few training samples, that is not surpassed by any other procedure that we know of, given the circumstances above. Based on the above, and on the fact that the K-L ordering is a very efficient technique in that it reduces the number of permutations of features that have to be searched through to only the number of features present, it will be used as a feature selection technique throughout the remainder of the experiments.

4.3 Experiments on the Hughes Phenomenon

In this section, some experimental results that illustrate the Hughes phenomenon will be presented. The objectives of conducting these experiments are t demonstrate the existence of this phenomenon in remote sensing applications, and to verify the hypothetical explanation of it. Experiments will be performed on aircraft and Landsat data, both simulated and real. In all the following experiments, no results are obtained for the dimensionality of one. Tabulated classification results are found in Appendix D.

Experiment 4.4

The data set described in experiment 4.1 is simulated using the algorithm described in section 2.4. Two classes, corn and forest, are selected and 500 training samples are chosen for each class. A larger, mutually exclusive set is

used for testing. The K-L method is used in ordering the features, and the data selected is classified using the best 2,3,4,..., 12 features. Subsequently, 5 training sets are randomly chosen from the larger training set, each set having 20 samples per class of corn and forest. The five sets are classified, using the same test fields above, and the average classification accuracy, (sometimes referred to as the probability of correct classification, or P_{cc}), is calculated for each subset of features. Another 5 training sets are then randomly chosen, this time with 13 samples per class of corn and forest (The minimum number of samples possible for 12 features without getting singular covariance matrices). Again, the 5 sets are classified and the average classification accuracy is calculated for each feature subset. The results are then plotted in Figure 4.4.

Looking at Figure 4.4, it is seen that when the number of training samples is adequate, as in the 500 samples per class case, the probability of correct classification is a monotonically non-decreasing function of dimensionality. Since in a K-L ordering, the information is concentrated in the first few channels, we notice that after the best 5 features, the recognition accuracy tends to saturate.

When the number of training samples per class drops to 20, however, we see that not only does the accuracy drop from the 500 samples case, but also it exhibits a slight Hughes phenomenon. Although the curve has a maximum at





dimensionality 3, it is approximately constant until the best 10 features, after which it starts decreasing, even though slightly.

The 13 samples per class case offers a dramatic change from the two other curves. There is a clear peaking effect here, with the curve reaching a maximum at dimensionality 5, after which it drops drastically.

The results conform with the hypothetical curves of Figures 2.1 and 2.2. The 20 samples and 13 samples curves can be made smoother if more than 5 sets are averaged, and hence we should look at them with the idea in mind that these are only approximations of what the true curves look like. However, the trend these curves point to is clear. In the presence of a limited set of training samples, an increase in dimensionality can result in a decrease in the classification accuracy, with this effect disappearing as the number of training samples increases.

Experiment 4.5

The same aircraft data set as that used in experiment 4.1 is used, but without any simulation. 400 samples each of corn and forest are selected for training, and a larger, separate set is used for testing. Again, 5 different training sets of 20 and 13 samples per class are randomly chosen from the original training set and classified. The average classification results for each feature subset are calculated and plotted. Results appear in Figure 4.5.

The curves in Figure 4.5 are not as smooth as they are in Figure 4.4. This is attributed to the fact that we are working with real data, which does not as well satisfy the assumptions we make as the simulated data does. Still, the curve with the 13 samples does generally poorer than the other two curves and drops dramatically in accuracy, whereas the 400 samples curve appears to saturate almost from the start. The 20 samples curve appears to have a slight peaking effect, although the curve is very noisy.

Experiment 4.6

The data set used in this experiment is obtained from Landsat, flown over Henry County, Indiana. To obtain a data set with more than the 4 features available from Landsat on any particular date, four data sets flown over the site at different times are used. The dates the data was collected on are: June 9, July 16, August 20, and September 26, all in 1978. The data is concatenated, and a K-L transformation was performed on it. Simulated data, more precisely meeting such assumptions as normality is generated, and the first 12 channels are used for classification. We will refer to this data as multitemporal data to indicate that it is collected over different times.



Figure 4.5 Experimental Classification Results of Aircraft, Real Data Using Different Numbers of Training Samples.

Two classes, corn and soybeans, are selected with 250 samples per class for training, and a larger independent s.t. for testing. Again, 5 different training sets of 20 and 13 samples per class are chosen from the original training set and classified. Results are averaged and plotted in Figure 4.6.

The same results obtained in the previous two experiments are again evident. Note that even with 20 or 13 samples per class, the accuracy obtained is very close to that obtained by using all the available training samples. This is due to the fact that the two classes chosen are highly separable and thus are easily distinguishible even when using a small number of training samples to estimate their statistics.

Experiment 4.7

The same data set as experiment 4.6 is used, but without any simulation. Two classes, corn and soybeans, are selected with 250 samples per class used for training, and a larger, separate, set for testing. Again, 5 different training sets of 20 and 13 samples per class are randomly chosen from the original training set and classified. Results are averaged and plotted in Figure 4.7.



Landsat, multitemporal, simulated data

Figure 4.6 Experimental Classification Results of Landsat, Multitemporal, Simulated Data Using Different Numbers of Training Samples.



Figure 4.7 Experimental Classification Results of Landsat, Multitemporal, Real Data Using Different Numbers of Training Samples.

The same observations noticed in the three previous experiments apply here. There is a drastic drop in accuracy when 13 samples are used, a slight one when cO samples are used, and no drop when 250 samples are used.

Summarizing the results of the last four experiments, we see that there is a definite Hughes phenomenon in the presence of a limited number of training samples compared to the number of features used. Further, as the number of samples increases, the accuracy for any given dimensionality increases, and the peak in the curve shifts to the right, i.e., the peaking effect takes place at a higher dimensionality, as is seen in Figures 4.4-4.7.

Studying Figures 4.4-4.7 reveals that the region between 13 samples and 20 samples is a very sensitive one when working with a maximum dimensionality of 12. While there is a sharp decline in accuracy at 13 samples per class, there is only a slight one at 20 samples per class. Another point to note is that the 20 and 13 samples are chosen from spectrally homogeneous classes, and so a very large number of samples is not needed to estimate the statistics of these classes. In a practical situation, the 20 and 13 samples curves might not be as close to the curves with large numbers of training samples as they are in these experiments.

The results of the last four experiments were a factor in choosing the empirical formula, or equation (3.61), discussed in Section 3.2.3. A formula was sought that takes the sensitivity in the number of training samples into account, as well as other factors that were discussed earlier.

4.4 Experiments Comparing Algorithm and Experimental Results

In this section, several experiments will be conducted to assess the performance of the proposed algorithm. Again, aircraft and Landsat data are used, both simulated and real, and the number of training samples used will be varied. But first, we will reproduce the results obtaired by Fukunaga and Krile (64) to verify the validity of the algorithm.

Experiment 4.8

The data set used by Fukunaga and Krile is described in detail in Marill and Green (12). The data is simulated, has two classes and eight features. Each class has 200 training samples, and both the exact, or true, and the algorithm recognition rates are calculated. The true recognition rates are not calculated again in here, but are reproduced from Fukunaga and Krile, who used numerical integration to arrive at them. Two methods used by Fukunaga and Krile are employed here: The normal assumption, discussed briefly in Section 3.2.1, and the modified gamma assumption, discussed in Section 3.2.2 and used throughout this research. The Bhattacharyya distance was used by Fukunaga and Krile, and although we have shown it to have limitations, it is used as a criterion for ordering the features. Results appear in Figure 4.8.

The results show that the modified gamma assumption method is a reasonable approximation of the true probability of correct classification. The normal assumption, though, does not give a good approximation of P_{CC} , and hence it is not further used.

While in this experiment, the modified gamma assumption is compared to the true probability of error, in actual practice the true probability of error cannot be calculated because the underlying distributions are not known. Therefore, in the following experiments, the proposed algorithm is compared to an average of five classifications obtained from five different training sets having the same number of training samples. This average classification serves as an estimate of the "true" error curve. This fact should be remembered as the experimental curves that are obtained are not as "smooth" as what the true curves would be expected to



Results Using Fukunaga and Krile

Figure 4.8 Classification Results of Fukunaga and Krile's Example Reproduced.

87 OF POOR QUALITY

be. The algorithm curves, on the other hand, being dependent, among other things, on the number of training samples in an average way, are expected to be "smoother" than the experimental ones.

Before we embark on studying the next experiments, it is appropriate at this point to look at a flowchart describing the modified algorithm that is proposed. This is shown in Figure 4.9. This figure is to be compared to Figure 3.2, or Fukunaga and Krile's algorithm, to see the changes that are made.

Experiment 4.9 (Aircraft, Simulated Data, 20 Samples per Class)

The simulated, aircraft data set used in Experiment 4.4 is used here. Two classes, corn and forest, are used. The experimental, 20 samples per class curve, in Figure 4.4 is plotted again in Figure 4.10, together with the approximation to the probability of correct classification predicted by the proposed algorithm. Also plotted in Figure 4.10 are the standard deviations for each feature subset of the five different classifications performed.

We see that the algorithm is a good approximation to the experimental curve. The approximation is not as good at lower dimensionalities as it is at higher ones, because the



Figure 4.9 A Flowchart of the Modified Algorithm.



Figure 4.10 Classification Results of Aircraft, Simulated Data, Using 20 Samples per Class.

C-2

assumptions the algorithm makes are better at higher dimensionalities. However, the two curves do peak at the same dimensionality, 3, but more importantly, they have a similar shape. Both remain relatively constant for a while and then start decreasing at about the dimensionality of 8.

Examining the standard deviations of P ..., it 19 observed that in general they have an increasing trend as the dimensionality increases. Put in other words, the curves indicate that the variance of the probability of error seems to increase with increasing dimensionality. This agrees with the hypothetical explanation given of the Hughes phenomenon, namely that the accuracy of the estimated statistics decreases with increasing dimensionality (i.e. becoming more random and hence increasing the variance of error) and that when this effect outweighs the increase in separability between classes due to increasing dimensionality, a peaking effect is observed. As the number of samples is decreased, larger increases in the variance of error are expected.

Experiment 4.10 (Aircraft, Simulated Data, 13 Samples per Class)

The same example used in Experiment 4.9 is used again, but with 13 samples per class used for training. The exper-

imental curve of Experiment 4.4 is reproduced, together with the curve predicted by the algorithm. The standard deviation of P_{cc} is again plotted. Results appear in Figure 4.11.

Again, the curve predicted by the algorithm is a better approximation of the experimental curve at high dimensionality. The experimental curve, however, is not very sensitive to dimensionality at lower values, and thus a small ambiguity in where the peak occurs can be afforded. Still, both curves predict a peak at 3. The standard deviation of the error again has an increasing trend as dimensionality increases.

Experiment 4.11 (Aircraft, Real Data, 20 Samples per Class)

The example used in Experiment 4.5 is repeated. Again, two classes are used, corn and forest, from the aircraft, real data set. Twenty samples per class are used for training, and five different sets of training samples are classified and averaged. The average is then compared to the algorithm performance. Results appear in Figure 4.12.

The experimental curve has a lot of error variance as can be seen from the curve and does not seem to be following any general pattern, although it starts consistently decreasing after dimensionality 9. It is interesting to



Aircraft, simulated data (13 samples/class)

Figure 4.11 Classification Results of Aircraft, Simulated Data, Using 13 Samples per Class.



Figure 4.12 Classification Results of Aircraft, Real Data, Using 20 Samples per Class.

compare this curve with Figure 4.10, where the same conditions exist with the exception that the data is simulated. Because simulated data satisfies the assumptions made about the distributions of classes, it produces results that conform more with theory than real data does. The algorithm performance appears to be closer to what is expected, although in this case it does not quite follow the experimental curve. This "randomness" of the experimental curve is made more evident from looking at the standard deviations of Pcc, which do not seem to follow any general pattern and are all relatively large. This is a clear example of a case where deviations from the assumptions may obscure the action of a new proposed algorithm.

Experiment 4.12 (Aircraft, Real Data, 13 Samples per Class)

The same example used in Experiment 4.11 is used here, with 13 samples per class for training. Results are shown in Figure 4.13.

Experimental and algorithm results here are very close. Both peak at 3, and both are very close at high dimensionalities. The standard deviations of the errors are also increasing in general, particularly at high dimensionality. It is interesting to note that the standard deviation in almost all of the above four experiments starts increasing

CARLES STATES



Figure 4.13 Classification Results of Aircraft, Real Data, Using 13 Samples per Class. notably at about the same place the probability of correct classification starts dropping sharply. This supports the idea that at these dimensionalities, the randomness in the estimated statistics is so large that it pulls the curve down.

Experiment 4.13 (Landsat, Multitemporal, Simulated Data, 20 Samples per Class)

The data set used in this experiment is the same as that used in Experiment 4.6. It is obtained from Landsat, with four dates concatenated so that more features are presented. The 20 samples per class curve of Figure 4.6 is reproduced in Figure 4.14, together with the curve predicted by the algorithm.

The algorithm curve seems to drop in accuracy faster than the experimental curve, but both peak at around 4. The standard deviation of error also increases as more features are used.

Experiment 4.14 (Landsat, Multitemporal, Simulated Data, 13 Samples per Class)

The same data set used in Experiment 4.13 is used, but with 13 samples per class for training. Results appear in



Landsat, multitemporal, simulated data (20 samples/class)

Figure 4.14 Classification Results of Landsat, Multitemporal, Simulated Data, Using 20 Samples per Class.
Figure 4.15. The increase in the variance of error with increasing dimensionality is very noticeable here. Again, the same observations apply, with both curves starting to drop in accuracy at the dimensionality of 4.

Experiment 4.15 (Landsat, Multitemporal, Real Data, 20 Samples per Class)

The Landsat data set is again used, but without any simulation. 20 samples per class are used for training, classification results are averaged and plotted in Figure 4.16.

While the algorithm predicts a somewhat better performance than the experimental curve, both have the same shape, and both are fairly constant until the first 7 or 8 features. This is due to the fact that the two classes in this set, corn and soybeans, are largely separable and hence the increase in the variance of the error with increasing dimensionality does not outweigh the large separability effect between these two classes.



Figure 4.15 Classification Results of Landsat, Multitemporal, Simulated Data, Using 13 Samples per Class.

99





Figure 4.16 Classification Results of Landsat, Multitemporal, Real Data, Using 20 Samples per Class.

Experiment 4.16 (Landsat, Multitemporal, Real Data, 13 Samples per Class)

The Landsat, real data set is used in this experiment with 13 samples per class for training. Results are shown in Figure 4.17. The two curves have the same shape, and peak at the same place, 4, although again the algorithm predicts a better performance than does the experimental curve. The variance of error is again seen to be increasing with the number of features used.

To summarize the results of the last eight experiments (4.9-4.16), the probabilities of error predicted by the proposed algorithm as a function of dimensionality as compared to experimental observations are shown for aircraft and Landsat data. Results are obtained for both simulated and real data, using 20 and 13 samples per class for training. For each case, five different training sets are used, and classification results are averaged over these five sets. The standard deviations of errors for each feature subset are also plotted.

Results indicate that the algorithm predicts in most of the cases the best, or near best, subset of features to be used. While not always predicting closely the actual clas-



Figure 4.17 Classification Results of Landsat, Multitemporal, Real Data, Using 13 Samples per Class.

sification accuracies obtained from the experimental average curve, it has in most of the cases the same shape as the experimental curve and seems to follow any trends in performance that the experimental curve undertakes. Since the objective behind the algorithm is to predict the best feature dimensionality and specific subset to be used in classification rather than to predict the probability of error itself, the fact that the algorithm does not always accurately predict this probability of error is not of serious concern.

The standard deviations plotted seem to indicate that in general, an increase in dimensionality results in an increase in the variance of error, that increase becoming highly noticeable at high dimensionality, when the randomness in the estimated statistics, given a limited set of training samples, is large.

The next step is to incorporate this algorithm in a binary tree classification procedure, using more than two classes, and assess its performance. This is done in Section 4.5.

4.5 Experiments on a Binary Tree Classification Procedure

In this section, two data sets will be classified in a binary tree classification procedure, using the proposed algorithm to predict the optimal features at every node.

A complete design of a binary tree classification procedure should address the problem of how to separate the nodes in the tree effectively. Seprations should be sought that lead to meaningful classes at the intermediate and terminal nodes. This problem should be thoroughly studied before a solution can be arrived at.

It is not the purpose of this research to address this problem in any detail. Therefore, no attempt has been made here to dictate a particular procedure or claim any optimal, or close to optimal, one. The procedure that will be used is heuristic, the purpose of conducting the next two experiments is to illustrate the usefuleness of the proposed algorithm in predicting the optimal features to be used at every node. The problem of how to separate the nodes is left as a topic for future research.

Experiment 4.17

The Landsat, multitemporal, real data set used in Experiment 4.6 is used here again. Three informational

classes exist in the scene: corn, soybeans, and other. 13 samples per class are used for training, creating 3 spectral classes. The reason this is done is that in actual practice situations, it is almost impossible to distinguish spectral classes with only 13 training samples per class. A much larger, separate, set is used for testing (all training and test field descriptions are found in Appendix F). The binary tree is constructed by using a bottom-up procedure, combining the most separable classes. The criterion for measuring separability is that used by Whitsict (9), and is defined as follows:

$$D_{erf} = erf((2B)^{1/2})$$
 (4.3)

where B is the Bhattacharyya distance and erf (.) is the gaussian error function. Whitsitt found that this measure is less ambiguous and more linear than the measure B. The measure is calculated using the first 12 features after a Karhunen-Loeve expansion was performed on the data. After the tree is constructed this way, the proposed algorithm was used to predict the optimal features to be used at every node.

The binary tree that resulted from the above procedure is shown in Figure 4.18. The algorithm predicts an optimal feature subset of 4 at the top, and a subset of 2 at the intermediate node. These appear below each node. Inside the node, the classes present are shown together with the total number of training samples present.



Figure 4.18 Binary Tree Design Structure of Landsat, Multitemporal, Real Data, Using 13 Training Samples per Class, With Numbers Inside Nodes Indicating Number of Training Samples Used.

AND AND MEDICAL STREET, SALAR STREET, SALAR STREET, SALAR STREET, SALAR STREET, SALAR STREET, SALAR STREET, SA

A single-stage classification is then performed on the data using feature subsets of 2 to 12. This is done to compare the performance of the binary tree procedure to bhat of each of the feature subsets.

Results are plotted in Figure 4.19. The classification result obtained from the binary tree procedure is drawn in a dotted line across the page only to compare against the single-stage curve, and does not imply that all the feature subsets were used, or that the classification result is the same for all feature subsets.

The results indicate that using three classes, the single-stage curve has a peak at 4, and that by using all twelve features, the result is much poorer. The binary tree procedure, on the other hand, results in a classification accuracy that is almost as good as the best result obtained from using the best feature subset (which is unknown in an actual practice situation) in a single-stage classification. Thus, it appears that the algorithm is effective in predicting the best features to be used at each node.

Experiment 4.18

The aircraft, real data set used in Experiment 4.1 is used here. The data set has seven informational classes.



Figure 4.19 Single-Stage and Binary Tree Classification Results of Landsat, Multitemporal, Real Data, Using 13 Training Samples per Class.

In this experiment, supervised clustering (discussed in Section 1.2.1) is used to get 9 spectral classes, using an adequate number of training samples per class. 13 samples per class were then randomly chosen from the larger training set so that it is known that each set of these samples comes from one spectral class. The bottom-up procedure described in Experiment 4.17 was then used to build the binary tree, with the exception of class water, which was separated from the other classes at the beginning, as water has been known from experience to have spectral properties that are much different from other agricultural classes. The proposed algorithm is then used to predict the best features at each node. A single-stage classification is performed using feature subsets of 2 to 12, and then the same statistics were used in the binary tree classification procedure.

The resulting tree appears in Figure 4.20. Figure 4.21 shows the classification results obtained from the single-stage and the binary tree classifiers.

The binary tree procedure, using the proposed algorithm, performs better than any feature subset does in a single-stage procedure. The Hughes phenomenom is very evident here, as the overall classification accuracy for seven informational classes (9 spectral) drops sharply from a high of 69.4% to a low of 43.0%.



· .

Figure 4.20 Binary Tree Design Structure of Aircraft, Real Data, Using 13 Training Samples per Class.



Figure 4.21 Single-Stage and Binary Tree Classification Results of Aircraft, Real Data, Using 13 Training Samples per Class.

Carlot Transfer 1

Summarizing the results of the last two experiments, the proposed algorithm is shown to be effective in predicting feature subsets that lead to the maximum, or near maximum, accuracy possible using the Karhunen-Loeve expansion for ordering the features.

It is worthwhile to note that common belief is that few features need be used at the top of the tree to separate classes, and more features need be used deeper in the tree to distinguish between somewhat inseparable classes. However, if there are inadequate training samples present, then the number of training samples towards the bottom of the tree is less than that towards the tor. Hence, less features should be used at the bottom to avoid the Hughes phenomenon. This is evident in the last two examples, particularly in Figure 4.20, where many features are used at the top, but only few at the bottom.

One point also worth mentioning is that in situations where a node is divided into two nodes of unequal training samples, one of them might have inadequate training samples while the other might have adequate ones. This situation is illustrated in Figure 4.20, where the top node is divided into water, and everything else. In this case, the number of features used is "intermediate", depending on the effect of the degradation in the accuracy of the estimated statistics of the node with the inadequate number of training samples.

CHAPTER 5

SUMMARY AND CONCLUSIONS

5.1 Summary of Results

The purpose of this research has been to develop an error estimator that will predict when/if the Hughes phenomenon occurs in multispectral data. Several significant results were arrived at and are summarized below.

The probability of error was studied through the likelihood ratio function, which offered the convenience of working with a one- dimensional variable, regardless of the number of features used in estimating the training statistics. An algorithm was then developed to estimate the statistics of this function, taking into account the number of training samples used to estimate these statistics. Several theoretical and experimental results were obtained on the Hughes phenomenon. These showed the dependency of the probability of error on the number of training samples and features used. The algorithm developed in Chapter 3 was shown to predict a suitable feature subset to be used at each node in a binary tree procedure. The algorithm was tested in Chapter 4 by comparing it to experimental observations under different conditions, and was utilized in two binary tree classification procedures to demonstrate its practicality.

Some results were also shown, demonstrating the usefuleness of the K-L expansion over the whole data set in ordering features in the presence of a limited set of training samples. The procedure is used extensively in the research, and appears to have less variablity than other procedures under the conditions given.

Certain parts of the algorithm developed are heuristic in nature. Reasons why more theoretical solutions were not pursued were explained. These heuristic procedures often raise difficulty in verifying the validity of the algorithm strategy. The basic point is that when both a practical solution and theoretical perfection cannot be achieved simultaneously, one tends tr choose the former. Experimental results in Chapter 4 demonstrated that the algorithm can be used practically to yield optimal, or near optimal, results.

5.2 Suggestions for Further Research

an and a state of the state of the second state of the second state of the state of the second state of th

The main objective behind developing the error algorithm is to use it as a feature selection technique in a multi-stage classification procedure. In particular, the algorithm was developed to be used in a binary tree procedure. The design of such a procedure requires, in addition

to choosing the optimal features at each node, an effective design of separating the nodes. This question was only addressed superficially in this research, and could serve as a topic for another research project. An effective design for separating the nodes, coupled with the developed algorithm to choose the features, should lead to much higher accuracies than a single-stage classifier.

Several strategies developed in the research were heuristic in nature. Appendix B addresses the problem of why it is difficult to theoretically calculate the probability density function of the variances of the likelihood ratio function given either class one or two. If such a derivation is made possible, a much better and clearer idea will be obtained on how the variance of the likelihood ratio function is affected by the number of training samples, and the error algorithm can be made to more accurately predict the probability of error in the presence of a limited number of training samples.

The K-L expansion was used extensively as a feature selection technique in the presence of few training samples. This was based on experimental observations, but necessarily meant sacrificing the information found from the between classes variablity. A more detailed study of the relation of several feature selection techniques to the number of training samples can be very helpful.

LIST OF REFERENCES

LIST OF REFERENCES

- Swain, P.H. and H. Hauska. 1977. The Decision Tree Classifier: Design and Potential. IEEE Trans. Geos. Elect. GE-15(3): 142-147.
- 2. Fukunaga, K. 1972. Introduction to Statistical Pattern Recognition. Academic Press, New York.
- 3. Duda, R.O. and P.E. Hart. 1973. Pattern Classification and Scene Analysis. Wiley, New York.
- 4. Fleming, M.D., J.S. Berkebile, and R.M. Hoffer. 1975. Computer-Aided Analysis of Landsat-1 MSS Data: A Comparison of Three Approaches, Including & 'Modified Clustering' Approach. 10p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Information Note 072475.
- 5. Fleming, M.D. and R.M. Hoffer. 1977. Computer-Aided Analysis Techniques for an Operational System to Map Forest Lands Utilizing Landsat MSS Data. 254p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Information Note 112277.
- 6. Fukunaga, K. and D.L. Kessell. 1973. Nonparametric Bayes Error Estimation Using Unclassified Samples. IEEE Trans. Infor. Theory. IT-19(7):434-400.
- Mobasseri, B.G. and C.D. McGillem. 1979. Multiclass Bayes Error Estimation by a Feature Space Sampling Technique. IEEE Trans. Systems, Man and Cybernetics. SMC-9(10):660-665.
- Moore, D.S., S.J. Whitsitt, and D.A. Landgrebe. 1976. Variance Comparisons of Unbiased Estimators of Probability of Correct Classification. IEEE Trans. Infor. Theory. IT-22(1):102-105.
- 9. Writsitt, S.J. and D.A. Landgrebe. 1977. Error Estimation and Separability Measures in Feature Selection for Multiclass Pattern Recognition. 186p. Laboratory for Applications of Remote Sensing, Purdue University, West

Lafayette, Indiana. LARS Publication 082377. Also available as a Ph.D Thesis, TR-EE 77-34. Department of Electrical Engineering, Purdue University. ħ

- 10. Wiersma, D.J. and D.A. Landgrebe. 1978. The Analytical Design of Spectral Measurements for Multispectral Remote Sensor Systems. 271p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Technical Report 122678. Also available as a Ph.D Thesis, TR-EE 79-13. Department of Elecrtical Engineering, Purdue University.
- 11. Mobasseri, B.G., D.J. Wiersma, E.R. Wiswell, D.A. Landgrebe, C.D. McGillem, and P.E. Anuta. 1978. A Multispectral Scanner System Parameter Study and Analysis Software System Description. 126p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Contract Report 112678.
- 12. Marill, T. and D.M. Green. 1963. On the Effectiveness of Receptors in Recognition Systems. IEEE Trans. Infor. Theory. IT-9(1):11-17.
- Swain, P.H. and R.C. King. 1973. Two Effective Feature Selection Criteria for Multispectral Remote Sensing.
 5p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Information Note 042673.
- 14. Jeffreys, H. 1948. Theory of Probability. Oxford University Press.
- Matusita, K. 1951. On the Theory of Statistical Decision Functions. Ann. Instit. Stat. Math. (Tokyo). 3:17-35.
- 16. Bhattacharyya, A. 1943. On a Measure of Divergence Between Two Statistical Populations Defined by Their Probability Distributions. Bul. Calcutta Math. Soc. 35:99-109.
- Mahalanobis, P.C. 1936. On the Generalized Distance in Statistics. Proc. National Inst. Sci. (India). 12:49-55.
- 18. Swain, P.H., T.V. Robertson, and A.G. Wacker. 1971. Comparison of the Divergence and B-Distance in Feature Selection. 12p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Information Note 020871.
- 19. Kailath, T. 1967. The Divergence and Bhattacharyya Distance Measures in Signal Selection. IEEE Trans. Communication Technology. COM-14(1):52-60.

- 20. Karhunen, K. 1947. Uber Lineare Methoden in Der Wahrscheinlich- keitsrechnung. Amer. Acad. Sci., Fennicade. Ser. A,I, 37:3-79. (Transl: Rand Corp., Santa Monica, California, Rept. T-131, August 1960.)
- 21. Loeve, M. 1963. Probability Theory. Van Nostrand, Princeton, New Jersey.
- 22. Oja, E. and J. Karhunen 1980. Recursive Construction of Karhunen-Loeve Expansions For Pattern Recognition Purposes. Proceedings of Fifth International Conference on Pattern Recognition, Miami Beach, Florida, December 1-4. pp. 1215-1218.
- 23. Karhunen, J. and E. Oja 1980. Some Comments on the Subspace Methods of Classification. Proceedings of Fifth International Conference On Pattern Recognition, Miami Beach, Florida, December 1-4. pp. 1191-1194.
- 24. Kanal, L. 1976. Patterns in Pattern Recognition: 1968-1974. IEEE Trans. Infor. Theory. IT-20(6):697-722.
- 25. Hughes, G.F. 1968. On the Mean Accuracy of Statistical Pattern Recognizer. IEEE Trans. Infor. Theory. IT-14(1):55-63.
- 26. Abend, K. and T.J. Harley, Jr. 1969. Comments "On the Mean Accuracy of Statistical Pattern Recognizers." IEEE Trans. Infor. Theory (Correspondence). IT(5):420-421.
- 27. Chandrasekaran, B. and T.J. Harley, Jr. 1969. Comments "On the Mean Accuracy of Statistical Pattern Recognizer." IEEE Trans. Infor. Theory (Correspondence). IT(5):421-423.
- 28. Kanal, L. and B. Chandrasekaran. 1971. On Dimensionality and Sample Size in Statistical Pattern Classification. Pattern Recognition. 3:225-236.
- 29. Chandrasekaran, B. 1971. Independence of Measurements and the Mean Recognition Accuracy. IEEE Trans. Information Theory, Vol IT-7(4):452-456.
- 30. Foley, D.H. 1972. Considerations of Sample and Feature Size. IEEE Trans. Information Theory. IT-18(5):618-626.
- 31. Chandrasekaran, B. and A.K. Jain 1975. Independence, Measurement Complexity and Classification Performance. IEEE Trans. Systems, Man and Cybernatics. SMC-5(2):240-244

THE REPORT OF TH

- 32. Duin, R.P. 1977. Comments on "Independence, Measurement Complexity, and Classification Performance". IEEE Trans. Systems Man and Cybernatics. (Correspondence) SMC(7):559-560.
- 33. Van Ness, J. 1977. Dimensionality and Classification Performance with Independent Coordinates. IEEE Trans. on Systems, Man and Cybernatics. (Correspondence) SMC(7):560-564.
- 34. Chandrasekaran, B. and A.K. Jain 1977. "Independence, Measurement Complexity and Classification Performance": An Ammendation. IEEE Trans. Systems, Man and Cybernatics. (Corrspondence) SMC(7):564-566.
- 35. Van Campenhout, J.M. 1978. On the Peaking of the Hughes Mean Recognition Accuracy: The Resolution of an Apparent Paradox. IEEE Trans. Systems, Man and Cybernetics. SMC-8(5):390-395.
- 36. Kulkarni, A.V. 1978. On the Mean Accuracy of Hierarchical Classifiers. IEEE Trans. Computers. C-27(8):771-776.
- 37. Raudys, S.J. 1979. Determination of Optimal Dimensionality in Statistical Pattern Classification. Pattern Recognition. 11:263-270.
- 38. Trunk, G.V. 1979. A Problem of Dimensionality: A Simple Example. IEEE Trans. Pattern Analysis and Machine Intelligence. PAMI-1:306-307.
- 39. Raudys, S. and V. Pikelis. 1980. On Dimensionality, Sample Size, Classification Error, and Complexity of Classification Algorithm in Pattern Recognition. IEEE Trans. Pattern Analysis and Machine Intelligence. PAMI-2(3):242-252.
- 40. El-Sheikh, T.S. and A.G. Wacker. 1980. Effect of Dimensionality and Estimation on the Performance of Gaussian Classifiers. Pattern Recognition 12:115-126.
- 41. Wacker, A.G. and T.S. El-Sheikh 1980. Calculation of Probability of Correct Classification for Two-Class Gaussian Classifiers With Arbitrary Hyperquadratic Decision Boundaries. Machine Processing of Remotely Sensed Data Symposium. CHI533-9:94-302.
- 42. El-Sheikh, T.S. and A.G. Wacker 1980. Average Classification Accuracy Over Collections of Gaussian Problems. CHI1499-3:685-690.
- 43. Wald, A. 1947. Sequential Analysis. Wiley, New York.

A CONTRACTOR OF A CONTRACTOR OF

44. Fu, K.S., Y.T. Chien and G.P. Cardillo. 1967. A Dynamic Programming Approach to Sequential Pattern Recognition. IEEE Trans. Computers. C-16(6):790-803.

÷

- 45. Fu, K.S. 1968. Sequential Methods in Pattern Recognition and Machine Learning. Academic Press.
- 46. Dubes, R. and A.K. Jain. 1979. Validity Studies in Clustering Methodologies. Pattern Recognition. 11:235-254.
- 47. Lukasova, A. 1979. Hierarchical Agglomerative Clustering Procedure. Pattern Recognition. 11:365-381.
- 48. Nadler, M. 1971. Error and Reject Rates in a Hierarchical Pattern Recognizer. IEEE Trans. Computers. C-20:1598-1601.
- 49. Meisel, W.S. and D.A. Michalopoulos. 1973. A Partitioning Algorithm with Application in Pattern Classification and the Optimization of Decision Trees. IEEE Trans. Computers. C-22:93-103.
- 50. Wu, C.L., D.A. Landgrebe and P.H. Swain. 1974. The Decision Tree Approach to Classification. 194p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Information Note 090174. Also available as a Ph.D Thesis, TR-EE 75-17. Department of Electrical Engineering, Purdue University.
- 51. Swain, P.H., C.L. Wu, D.A. Landgrebe, and H. Hauska. 1975. Layered Classification Techniques for Remote Sensing Applications. 12p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Information Note 061275.
- 52. Bartolucci, L.A., P.H. Swain, and C.L. Wu. 1976. Selective Radiant Temperature Mapping Using a Layered Classifier. IEEE Trans. Geoscience Electronics. GE-14:101-106.
- 53. You, K.C. and K.S. Fu. 1976. An Approach to the Design of a Linear Binary Tree Classifier. Proc. Conference on Machine Processing of Remotely Sensed Data. June 29-July 1, 1976. IEEE Catalog No. 76CH1103-1 MPRSD.
- 54. Fletcher, R. and M.J.D. Powell. 1963. A Rapid Descent Method for Minimization. Computer Journal. 6:163-168.
- 55. Kulkarni, A.V. and L.N. Kanal. 1976. An Optimization Approach to Hierarchical Classifier Design. Proc. Third Int. Joint Conf. on Pattern Recognition (Coronado, California), IEEE Catalog No. 76CH/140-3C.

- 56. Parikh, J. 1977. A Comparative Study of Cloud Classification Techniques. Remote Sensing of Environment. 6:67-81.
- 57. Sethi, I.K. and B. Chatterjee. 1977. Efficient Decision Tree Design for Discrete Variable Pattern Recognition Problems. Pattern Recognition. 9:197-206.
- 58. Breiman, L. 1978. Growing Trees to Analyze High Dimensional Data. Technology Service Corporation Report. TSC-CSD-IN-024.
- 59. Sonquist, J.A., E.L. Baker, and J.N. Morgan. 1973. Searching for Structure. Survey Research Center, Institute for Social Research, University of Michigan, Ann Arbor, Michigan.
- 60. Rounds, E.M. 1979. A Combined Nonparametrical Approach to Feature Selection and Binary Decision Tree Design. Proc. IEEE Computer Society Conference on Pattern Recognition and Image Processing, Chicago, Illinois, August 6-8. CH1428-2:38-43.
- 61. Dattatreya, G.R., and V.V.S. Sarma, 1981. Bayesian and Decision Tree Approaches for Pattern Recognition Including Feature Measurements Costs. IEEE Trans. Patvern Recognition and Machine Intelligence. PAMI-3:293-298.
- 62. Wacker, A.G. and D.A. Landgrebe. 1971. The Minimum Distance Approach to Classification. 361p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Information Note 100771. Also available as a Ph.D Thesis, TR-EE 71-37, Department of Electrical Engineering, Purdue University.
- 63. Kullback, S. 1959. Information Theory and Statistics. p. 195. Wiley, New York.
- 64. Fukunaga, K. and T. Krile. 1969. Calculation of Bayes Recognition Error for Two Multivariate Gaussian Distributions. IEEE Trans. on Computers. C-18(3):220-229.
- 65. Swain, P.H. and S.M. Davis, eds. 1978. Remote Sensing: The Quantitative Approach. pp. 164-174. McGraw-Hill, Inc., New York.
- 66. Muasher, M. and P. Swain. 1980. A Multispectral Data Simulation Technique. 30p. Laboratory for Applications of Remote Sensing, Purdue University, West Lafayette, Indiana. LARS Technical Report 070980.

- 67. Van Trees, H.L. 1968. Detection, Estimation and Modulation Theory, Part I. Wiley & Sons. New York.
- 68. Bickel, P.J. and K.A. Doksum 1977. Mathematical Statistics: Basic Ideas and Selected Topics. Holden-Day, San Francisco.
- 69. Papoulis, A. 1965. Probability, Random Variables, and Stochastic Processes. McGraw-Hill, New York.

and the subsequences

APPENDICES

ž

Appendix A

Generation of Normally Distributed Samples

Let U_1 and U_2 be two random variables independent and identically distributed Uniform (0,1). Then, let

$$Z_{1} = (-2 \ln U_{1})^{\frac{1}{2}} \cos 2\pi U_{2}$$
 (A.1)

$$Z_{2} = (-2 \ln U_{1})^{\frac{1}{2}} \sin 2\pi U_{2}$$
 (A.2)

then Z_1 and Z_2 are independent and identically distributed normal (0, 1).

Proof:

$$f(U_1, U_2) = \begin{cases} 1 & 0 < U_1 < 1, 0 < U_2 < 1 \\ 0 & \text{otherwise} \end{cases}$$
(A.3)

is the probability density function of two independent uniforms.

$$U_{1} = \exp\left[-\frac{1}{2}(z_{1}^{2} + z_{2}^{2})\right]$$
(A.4)
$$U_{2} = \frac{1}{2\pi} \arctan\left(\frac{z_{2}}{z_{1}}\right)$$
(A.5)

The jacobian of the transformation is:

$$J = -\frac{1}{2\pi} \exp\left[-\frac{1}{2}(Z_{1}^{2} + Z_{2}^{2})\right]$$

$$f(Z_{1}, Z_{2}) = f(U_{1}, U_{2}) \cdot \left[J\right]$$

$$= \frac{1}{2\pi} \exp\left[-\frac{1}{2}(Z_{1}^{2} + Z_{2}^{2})\right] \quad 0 < \exp\left[-\frac{1}{2}(Z_{1}^{2} + Z_{2}^{2})\right] < 1$$

$$0 < \frac{1}{2\pi} \arctan\left(\frac{Z_{2}}{Z_{1}}\right) < 1$$

$$= 0 \quad \text{otherwise} \qquad (A.6)$$

$$f(Z_1) \sim N(0,1)$$
 $f(Z_2) \sim N(0,1)$

The side conditions give $-\infty < Z_1 < \infty$, $-\infty < Z_2 < \infty$. Strictly speaking, Z_1 cannot equal zero; however, prob $(Z_1 = 0)=0$ as we are working with continuous densities.

To test the effectiveness of the pseudo random vectors ir the multivariate case, random vectors distributed $N(0,I_p)$ were generated and then tested with a Kolmogorov-Smirnov test. Since the multivariate normal cdf is difficult to evaluate, the sum of squares was calculated and compared to the x_p^2 distribution.

For sample sizes greater than 100, the pseudo random vectors were distributed properly. For sample sizes less than 100, the K-S test is not valid. Since we would generally (over an entire area) be working with more than 100 points per class, this was not pursued further.

In addition, the sample covariance matrices were tested for homogeneity against the true class statistics. For sample runs of up to 2000 points, there were not significant differences at the $\alpha = 0.10$ level.

Appendix B

On The Probability Density Functions

Of $\hat{\sigma}_1^2$ And $\hat{\sigma}_2^2$

Let us look at the expressions for $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$. From (3.55) and (3.58), we have:

$$\hat{\sigma}_{1}^{2} = 2(tr(I - \hat{\Sigma}_{2}^{-1} \hat{\Sigma}_{1})^{2} + 2\hat{m}_{2} \hat{\Sigma}_{2}^{-1} \hat{\Sigma}_{1} \hat{\Sigma}_{2}^{-1} \hat{m}_{2}) \quad (B.1)$$

$$\hat{\sigma}_{2}^{2} = 2(tr(\hat{\Sigma}_{1}^{-1} \hat{\Sigma}_{2} - I)^{2} + 2\hat{m}_{2} \hat{\Sigma}_{1}^{-1} \hat{\Sigma}_{2} \hat{\Sigma}_{1}^{-1} \hat{m}_{2}) \quad (B.2)$$

To be able to calculate the probability density functions of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$, one has to know those of \hat{m}_2 , $\hat{\Sigma}_1$, $\hat{\Sigma}_1^{-1}$, $\hat{\Sigma}_2$, and $\hat{\Sigma}_2^{-1}$.

Before we proceed, we make the following assumptions: 1. \hat{M}_1 and \hat{M}_2 , the means of the two classes at hand are constant. Experience has shown that one can estimate these two quantities relatively accurately with a small number of training samples. Henceforth, we will assume \hat{m}_2 (= $\hat{M}_1 - \hat{M}_2$) to be constant and not a random variable.

ORIGINIAL PAGE IS OF POOR QUALITY.

2. Σ_{1} and Σ_{2} are independent. We will ignore any relationships that might exist between the covariance matrices of the two classes.

Theorem B.

 $\hat{\Sigma}_1$, $\hat{\Sigma}_2$ are each Wishart distributed with parameters $\frac{1}{n_1}\Sigma_1$, n_1 and $\frac{1}{n_2}\Sigma_2$, n_2 respectively, where $n_1 = N_1 - 1$ and N_1 is the number of samples used in estimating Σ_1 .

Proof

See (B.1), pp. 159.

Thus, $\hat{\Sigma}_{i}$, i=1,2, has the following Wishart distribution:

$$\hat{\Sigma}_{i} \sim \frac{\binom{n_{i}}{2} \frac{\hat{\Sigma}_{i}}{2}}{\frac{n_{i}}{2} \frac{p}{\pi}} \frac{\hat{\Sigma}_{i}}{p(p-1)/4} \frac{\frac{n_{i}-p-1}{2}}{\frac{n_{i}}{2} \frac{p}{\pi}} \frac{\exp(-\frac{1}{2}(n_{i} \operatorname{tr} \hat{\Sigma}_{i}^{-1} \hat{\Sigma}_{i}))}{\frac{n_{i}}{2} \frac{p}{\pi}} (B.3)$$

where p is the number of dimensions.

Theorem B.2

 $\hat{\Sigma}_{i}^{-1}$ is again Wishart distributed with parameters $\frac{1}{n}_{i}\hat{\Sigma}_{i}^{-1}$, n_{i} . Proof

See (B.2)

Theorem B.3

If A is distributed according to Wishart, W(Σ ,n), then B = CAC^T is also distributed Wishart W(Φ ,n), where 4 = $C \Sigma C^{T}$.

Proof

1

NULLER DIS CONTRACTOR

See (B.1),pp.162.

From the above theorems, we see that $\hat{\Sigma}_1$, $\hat{\Sigma}_2$, $\hat{\Sigma}_1^{-1}$, and $\hat{\Sigma}_2^{-1}$ are Wishart distributed. Further, as $\hat{\Sigma}_1$ is transformed into the identity matrix I, and $\hat{\Sigma}_2$ is transformed into a diagonal matrix Λ , the new covariance matrices are also Wishart distributed. Hence, $\hat{\Sigma}_1$ is transformed into a diagonal matrix $\hat{\Gamma}$ that is distributed W(1/n_1I,n_1). We will call the diagonal elements of this matrix $\hat{\gamma}_1$. Similarly, $\hat{\Sigma}_2$ is transformed into a diagonal matrix $\hat{\lambda}_1$. $\hat{\Sigma}_1^{-1}$ is transformed into a diagonal matrix $\hat{\chi}_1$. Similarly, $\hat{\Sigma}_2$ is transformed into a diagonal matrix $\hat{\chi}_1$. $\hat{\Sigma}_1^{-1}$ is transformed into a diagonal matrix $\hat{\Gamma}^{-1}$ distributed W(1/n_1I,n_1), and $\hat{\Sigma}_2^{-1}$ is transformed into a diagonal matrix $\hat{\Gamma}^{-1}$ distributed W(1/n_2 Λ ,n_2).

Thus, after applying the simultaneous diagonalization transformation, $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ become:

$$\hat{\sigma}_{1}^{2} = 2 \sum_{i=1}^{p} (1 - \frac{\hat{\gamma}_{i}}{\hat{\lambda}_{i}} + \frac{\hat{\gamma}_{1}^{2}}{\hat{\lambda}_{i}^{2}} + 2 d_{1}^{2} \frac{\hat{\gamma}_{i}}{\hat{\lambda}_{i}^{2}})$$
(B.4)

$$\hat{\sigma}_{2}^{2} = 2 \sum_{i=1}^{p} \left(\frac{\hat{\lambda}_{i}^{2}}{\hat{\gamma}_{i}^{2}} - 2 \frac{\hat{\lambda}_{i}}{\hat{\gamma}_{i}} + 2 d_{i}^{2} \frac{\hat{\lambda}_{i}}{\hat{\gamma}_{i}^{2}} + 1 \right)$$
(B.5)

Note that equations (B.4) and (B.5) are modified versions of equations (3.53) and (3.56).

We now look at the probability density functions of the one-dimensional elements $\hat{\lambda}_i$ and $\hat{\gamma}_i.$

Theorem B.4

If $\Sigma_{ij}=0$ for $i \neq j$, and if A is distributed according to W(Σ , n), then A₁₁, A₂₂, ..., A_{pp} are independently distributed and A_{jj} is distributed according to W(Σ_{jj} , n). <u>Proof</u>

See (B.1),pp.163.

Therefore, $\hat{\lambda}_1, \ldots, \hat{\lambda}_p$ are each distributed W($\frac{\lambda_i}{n_2}, n_2$) and $\hat{\gamma}_1, \ldots, \hat{\gamma}_p$ are each distributed W(1/nj,nj). Hence,

 $\hat{\gamma}_{i} = \sqrt{\begin{cases} \frac{\hat{\gamma}_{i}(n_{1}-2)/2}{\hat{\gamma}_{i}} & \exp(-\frac{1}{2}n_{1}\hat{\gamma}_{i}) & (n_{1}/2) \\ & \Gamma(n_{1}/2) \\ 0 & & \hat{\gamma}_{i} < 0 \end{cases}} & (B.6)$

A similar expression exists for $\hat{\gamma_i}^{-1}$, with $\hat{\gamma_i}$ replaced by $\hat{\gamma_i}^{-1}$.

$$\hat{\lambda}_{i} \sim \begin{cases} \frac{\hat{\lambda}_{i}^{(n_{2}-2)/2}}{\sum_{i} \frac{n_{2}}{2} \sum_{i} \frac$$

A similar expression exists for $\hat{\lambda}_{i}^{-1}$, with $\hat{\lambda}_{i}$, λ_{i} replaced by $\hat{\lambda}_{i}^{-1}$, λ_{i}^{-1} .

Looking at equations (B.5) and (B.6), we see that even though we know the individual distributions of $\hat{\lambda}_1$ and $\hat{\gamma}_1$, the calculation of the densities of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ is still a very involved and difficult process. An attempt to arrive at these densities directly from those expressions is almost impossible. However, the moments of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ can be calculated.

Since calculating the moments of $\hat{\lambda}_i$ (and $\hat{\lambda}_i^{-1}$, $\hat{\gamma}_i$, $\hat{\gamma}_i^{-1}$) involves the evaluation of an integral of the type $\int^{\infty} t^n e^{-at} dt$, and since such an integral does indeed exist, 0 the task of calculating any moment of $\hat{\lambda}_i$, $\hat{\lambda}_i^{-1}$, $\hat{\gamma}_i$, and $\hat{\gamma}_i^{-1}$ is a very easy one.

From any integration table book, we find:

$$\int_{0}^{\infty} t^{n} \exp(-at) dt = \frac{\Gamma(n+1)}{a^{n+1}} \quad (n \ge 1, a \ge 0) \quad (B.8)$$

Thus, if x is distributed W(x/n,n), then:

STRUCTURE CONTRACTOR OF A CONTRACT OF A CONT
$$E(\mathbf{x}) = \mathbf{x}$$

$$E(\mathbf{x}^2) = (1+2/n) \mathbf{x}^2$$

$$E(\mathbf{x}^3) = (1+6/n + 8/n^2) \mathbf{x}^3$$

$$E(\mathbf{x}^4) = (1+12/n + 44/n^2 + 48/n^3) \mathbf{x}^4$$
(B.9)

Since any moment of $\hat{\sigma}_1^2$ or $\hat{\sigma}_2^2$ is a function of the moments of $\hat{\lambda}_1$, $\hat{\lambda}_1^{-1}$, $\hat{\gamma}_1$, and $\hat{\gamma}_1^{-1}$, it is theoretically possible to calculate any moment of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$. Thus, it is theoretically possible to calculate the characteristic function of $\hat{\sigma}_1^2$ or $\hat{\sigma}_2^2$ uniquely from these moments.

Papoulis (B.3) provides a way to estimate the probability density function of a random variable once its characteristic function is known. However, the convergence properties of calculating the characteristic function from the moments of a random variable are very slow. A large number of moments would have to be calculated. Looking at equations (B.4) and (B.5), it is evident that beyond the first few moments, the derivation becomes quite a formidable task, and is very impractical.

Because of these difficulties encountered, it was decided to calculate only the variances of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ and heuristically incorporate them into the algorithm developed.

REFERENCES

- B.1 Anderson, T.W. 1958. Introduction to Multivariate Statistical Analysis. Wiley, New York.
- B.2 Keehn, D.G. 1965. A Note On Learning For Gaussian Properties. IEEE Trans. Information Theory. pp. 126-132.
- B.3 Papoulis, A. 1962. The Fourrier Integral And Its Applications. McGraw-Hill, New York.

NERENALISTIC AVANTIC - -----

Appendix C

Derivation of the Variances of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$

We look first at $\hat{\sigma}_1^2$

From Appendix B, equation (B.4), we have

$$\hat{\sigma}_{1}^{2} = 2 \sum_{i=1}^{p} \left[1 - 2 \frac{\hat{\gamma}_{i}}{\hat{\lambda}_{i}} + \frac{\hat{\gamma}_{i}^{2}}{\hat{\lambda}_{i}^{2}} + 2 d_{i}^{2} \frac{\hat{\gamma}_{i}^{2}}{\hat{\lambda}_{i}^{2}} \right]$$
(C.1)

Noting the assumption that the $\hat{\lambda}_i$'s are independent from the $\hat{\gamma}_i^2$, and taking the expected value of $\hat{\sigma}_1^2$, we get

$$E(\hat{\sigma}_{1}^{2}) = 2 \sum_{i=1}^{p} \left[1 - 2 \frac{E(\hat{\gamma}_{i})}{E(\hat{\lambda}_{i})} + \frac{E(\hat{\gamma}_{i}^{2})}{E(\hat{\lambda}_{i}^{2})} + 2 d_{i}^{2} \frac{E(\hat{\gamma}_{i}^{2})}{E(\hat{\lambda}_{i}^{2})} \right]$$
(C.2)

Making use of the expressions in (B.9), we get

$$E(\hat{\sigma}_{1}^{2}) = 2\sum_{i=1}^{p} \left[1 - \frac{2}{\lambda_{i}} + (1 + \frac{2}{n_{1}})(1 + \frac{2}{n_{2}})\frac{1}{\lambda_{i}^{2}} + 2d_{i}^{2}(1 + \frac{2}{n_{2}})\frac{1}{\lambda_{i}^{2}}\right](C.3)$$

Now note that $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ are the summation of uncorrelated random variables. Since $\hat{\lambda}_i$'s are independent, $\hat{\gamma}_i$'s are independent, and each $\hat{\lambda}_i$ is independent from each $\hat{\gamma}_i$, then any function of $\hat{\lambda}_i$'s and $\hat{\gamma}_i$'s in one dimension is uncorrelated with any other function of $\hat{\lambda}_i$'s and $\hat{\gamma}_i$'s in another dimension. Hence, the variances of $\hat{\sigma}_1^2$ and $\hat{\sigma}_2^2$ consist of the sum of the variances in each dimension (See (69), p. 211) and

do not have any cross-product terms between dimensions. Therefore, in the following derivations, we will not attempt to derive any cross-product terms as they will cancel out in the end result.

$$\begin{bmatrix} E(\hat{g}_{1}^{2})^{2} \end{bmatrix} = 4E \left(\sum_{i=1}^{p} \left(1 - 2 \frac{\hat{y}_{i}}{\hat{\lambda}_{i}} + \frac{\hat{y}_{i}^{2} + 2d_{i}^{2}\hat{y}_{i}}{\hat{\lambda}_{i}^{2}} \right) \right)^{2} + \frac{\text{cross-product}}{\text{terms}}$$

$$= 4E \sum_{i=1}^{p} \left(1 - 4 \frac{\hat{y}_{i}}{\hat{\lambda}_{i}} + 2 \left(\frac{\hat{y}_{i}^{2} + 2d_{i}^{2}\hat{y}_{i}}{\hat{\lambda}_{i}^{2}} \right) + 4 \frac{\hat{y}_{i}^{2}}{\hat{\lambda}_{i}^{2}} - 4 \frac{\hat{y}_{i}^{3}}{\hat{\lambda}_{i}^{3}} - \frac{\hat{y}_{i}^{3}}{\hat{\lambda}_{i}^{3}} - \frac{1}{\hat{\lambda}_{i}^{3}} + \frac{\hat{y}_{i}^{4} + 4d_{i}^{2}\hat{y}_{i}^{3} + 4d_{i}^{4}\hat{y}_{i}^{2}}{\hat{\lambda}_{i}^{4}} \right) + \frac{\text{cross-product}}{\text{terms}}$$

$$= 4E \sum_{i=1}^{p} \left(1 - 4 \frac{\hat{y}_{i}}{\hat{\lambda}_{i}} + 2 \left(\frac{\hat{y}_{i}^{2} + 2d_{i}^{2}\hat{y}_{i}}{\hat{\lambda}_{i}^{2}} \right) + 4 \frac{\hat{y}_{i}^{2}}{\hat{\lambda}_{i}^{2}} - 4 \frac{\hat{y}_{i}^{3}}{\hat{\lambda}_{i}^{3}} - \frac{1}{\hat{\lambda}_{i}^{3}} - \frac{1}{\hat{\lambda}_{i}^{3}} + \frac{\hat{y}_{i}^{4} + 4d_{i}^{2}\hat{y}_{i}^{3} + 4d_{i}^{4}\hat{y}_{i}^{2}}{\hat{\lambda}_{i}^{4}} + \frac{1}{\hat{\lambda}_{i}^{4}} + \frac$$

Substituting the expressions of (B.9) into (C.4), we get

$$\begin{bmatrix} E(\hat{\sigma}_{1}^{2})^{2} \end{bmatrix} = 4 \sum_{i=1}^{p} \left(1 - \frac{4}{\lambda_{i}} + 2 \left(\frac{(1+2/n_{1}) + 2d_{1}^{2}}{\lambda_{1}^{2}} \right) (1 + \frac{2}{n_{2}}) \right)$$

$$+ \frac{4}{\lambda_{1}^{2}} (1 + \frac{2}{n_{1}}) (1 + \frac{2}{n_{2}}) - \frac{4}{\lambda_{1}^{3}} (1 + \frac{6}{n_{1}} + \frac{8}{n_{1}^{2}}) (1 + \frac{6}{n_{2}} + \frac{8}{n_{2}^{2}})$$

$$- \frac{8d_{1}^{2}}{\lambda_{1}^{3}} (1 + \frac{2}{n_{1}}) (1 + \frac{6}{n_{2}} + \frac{8}{n_{2}^{2}}) + \frac{1}{\lambda_{1}^{4}} (1 + \frac{12}{n_{2}} + \frac{44}{n_{2}^{2}} + \frac{48}{n_{2}^{3}})$$

$$\left((1 + \frac{12}{n_{1}} + \frac{44}{n_{1}^{2}} + \frac{48}{n_{1}^{3}}) + 4d_{1}^{2} (1 + \frac{6}{n_{1}} + \frac{8}{n_{1}^{2}}) + 4d_{1}^{4} (1 + \frac{2}{n_{1}}) \right) \right)$$

$$+ \frac{\text{cross-product}}{\text{terms}}$$

All of Bank Surgers of the

$$= 4 \frac{p}{1-1} \left(1 - \frac{4}{\lambda_{1}} + \frac{2}{\lambda_{1}^{2}} \left(1 + \frac{2}{n_{1}} + \frac{2}{n_{2}} + \frac{4}{n_{1}n_{2}} + 2d_{1}^{2} + 4\frac{d_{1}^{2}}{n_{2}} \right) \right) \\ + \left(4 + \frac{8}{n_{1}} + \frac{8}{n_{2}} + \frac{16}{n_{1}n_{2}} \right) \frac{1}{\lambda_{1}^{2}} - \frac{4}{\lambda_{1}^{2}} \left(1 + \frac{6}{n_{1}} + \frac{6}{n_{2}} + \frac{8}{n_{1}^{2}} + \frac{8}{n_{2}^{2}} + \frac{36}{n_{1}n_{2}} \right) \\ + \frac{48}{n_{1}n_{2}^{2}} + \frac{48}{n_{1}^{2}n_{2}} + \frac{64}{n_{1}^{2}n_{2}^{2}} \right) - 8\frac{d_{1}^{2}}{\lambda_{1}^{4}} \left(1 + \frac{2}{n_{1}} + \frac{6}{n_{2}} + \frac{12}{n_{1}n_{2}} + \frac{8}{n_{2}^{2}} + \frac{36}{n_{1}n_{2}^{2}} \right) \\ + \frac{1}{\lambda_{1}^{4}} \left(1 + \frac{12}{n_{1}} + \frac{12}{n_{2}} + \frac{144}{n_{1}n_{2}} + \frac{44}{n_{1}^{2}} + \frac{44}{n_{2}^{2}} + \frac{48}{n_{1}^{3}} + \frac{48}{n_{2}^{3}} + \frac{528}{n_{1}^{2}n_{2}} + \frac{528}{n_{2}^{2}n_{1}} \right) \\ + \frac{1936}{n_{1}^{2}n_{2}^{2}} + \frac{576}{n_{1}^{2}n_{2}} + \frac{576}{n_{2}^{2}n_{1}} + \frac{2112}{n_{1}^{2}n_{2}^{2}} + \frac{2304}{n_{1}^{2}n_{2}^{2}} + 4d_{1}^{2} \left(1 + \frac{6}{n_{1}} + \frac{12}{n_{2}^{2}} \right) \\ + \frac{8}{n_{1}^{2}} + \frac{44}{n_{2}^{2}} + \frac{72}{n_{1}n_{2}} + \frac{264}{n_{1}n_{2}^{2}} + \frac{96}{n_{1}^{2}n_{2}^{2}} + \frac{48}{n_{1}^{2}n_{2}^{2}} + \frac{352}{n_{1}^{2}n_{2}^{2}} + \frac{384}{n_{1}^{2}n_{2}^{2}} \right) \\ + 4d_{1}^{4} \left(1 + \frac{2}{n_{1}} + \frac{12}{n_{2}} + \frac{264}{n_{1}^{2}n_{2}^{2}} + \frac{24}{n_{1}n_{2}} + \frac{48}{n_{2}^{3}} + \frac{288}{n_{1}^{2}n_{2}^{2}} + \frac{352}{n_{1}^{2}n_{2}^{2}} + \frac{384}{n_{1}^{2}n_{2}^{2}} \right) \right) \\ + 4d_{1}^{4} \left(1 + \frac{2}{n_{1}} + \frac{12}{n_{2}} + \frac{44}{n_{2}^{2}} + \frac{24}{n_{1}^{2}n_{2}} + \frac{48}{n_{1}^{2}} + \frac{288}{n_{1}^{2}n_{2}^{2}} + \frac{352}{n_{1}^{2}n_{2}^{2}} + \frac{384}{n_{1}^{2}n_{2}^{2}} \right) \right) \\ + 4d_{1}^{4} \left(1 + \frac{2}{n_{1}} + \frac{12}{n_{2}} + \frac{44}{n_{2}^{2}} + \frac{24}{n_{1}^{2}n_{2}} + \frac{48}{n_{1}^{2}} + \frac{288}{n_{1}^{2}n_{2}^{2}} + \frac{352}{n_{1}^{2}n_{2}^{2}} + \frac{384}{n_{1}^{2}n_{2}^{2}} \right) \right) \right) \\ + \frac{2d_{1}^{2}}{n_{1}} \left(1 - \frac{2}{n_{1}} + \frac{12}{n_{2}} + \frac{24}{n_{1}^{2}n_{1}} + \frac{2}{n_{2}^{2}} + \frac{36}{n_{1}^{2}n_{2}^{2}} \right) \frac{1}{n_{1}^{2}}} \right) \\ + \frac{2d_{1}^{2}}{n_{1}^{2}} \left(1 + \frac{2}{n_{2}} \right) \right)^{2} + \frac{cross - product}{terms}$$

$$-\frac{4}{\lambda_{1}^{3}}\left(1+\frac{2}{n_{1}}+\frac{2}{n_{2}}+\frac{4}{n_{1}n_{2}}+2d_{1}^{2}\left(1+\frac{2}{n_{2}}\right)\right)+\frac{1}{\lambda_{1}^{4}}\left(1+\frac{4}{n_{1}}+\frac{4}{n_{2}}\right)$$

$$+\frac{4}{n_{2}^{2}}+\frac{4}{n_{2}^{2}}+\frac{16}{n_{1}n_{2}}+\frac{16}{n_{1}^{2}n_{2}}+\frac{16}{n_{1}n_{2}^{2}}+\frac{16}{n_{1}^{2}n_{2}^{2}}+4d_{1}^{2}\left(1+\frac{2}{n_{1}}+\frac{4}{n_{2}}\right)$$

$$+\frac{8}{n_{1}n_{2}}+\frac{4}{n_{2}^{2}}+\frac{8}{n_{1}n_{2}^{2}}\right)+4d_{1}^{4}\left(1+\frac{4}{n_{2}}+\frac{4}{n_{2}^{2}}\right)\right)+\frac{\operatorname{cross-product}}{\operatorname{terms}}$$

$$(C.6)$$
Now way $\left(\frac{2}{n_{1}}\right) = \int \mathbb{P}\left(\frac{2}{n_{1}^{2}}\right)^{2} = \int \mathbb{P}\left(\frac{2}{n_{1}^{2}}\right)^{2}$

Now, Var $(\hat{\sigma}_1^2) = [E(\hat{\sigma}_1^2)^2] - [E(\hat{\sigma}_1^2)]^2$ or,

1900 C

$$\begin{aligned} \operatorname{Var}\left(\hat{\mathfrak{o}}_{1}^{2}\right) &= 4 \Pr_{1=1}^{P} \left(\frac{2}{\lambda_{1}^{2}} \left(\frac{4}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{1}n_{2}} \right) - \frac{4}{\lambda_{1}^{3}} \left(\frac{4}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{1}^{2}} \right) \\ &+ \frac{8}{n_{2}^{2}} + \frac{32}{n_{1}n_{2}} + \frac{48}{n_{1}n_{2}^{2}} + \frac{48}{n_{1}^{2}n_{2}} + \frac{64}{n_{1}^{2}n_{2}^{2}} + \frac{4d_{1}^{2}}{n_{1}} + \frac{8d_{1}^{2}}{n_{2}} + \frac{24d_{1}^{2}}{n_{1}n_{2}} + \frac{16d_{1}^{2}}{n_{2}^{2}} \\ &+ \frac{32d_{1}^{2}}{n_{1}n_{2}^{2}} \right) + \frac{1}{\lambda_{1}^{4}} \left(\frac{8}{n_{1}} + \frac{8}{n_{2}} + \frac{128}{n_{1}n_{2}} + \frac{40}{n_{1}^{2}} + \frac{40}{n_{2}^{2}} + \frac{48}{n_{1}^{3}} + \frac{48}{n_{2}^{3}} + \frac{512}{n_{1}^{2}n_{2}} \\ &+ \frac{1920}{n_{1}^{4}n_{2}^{2}} + \frac{576}{n_{1}^{3}n_{2}} + \frac{576}{n_{2}^{3}n_{1}} + \frac{2112}{n_{1}^{2}n_{2}^{3}} + \frac{2112}{n_{1}^{3}n_{2}^{2}} + \frac{2304}{n_{1}^{3}n_{2}^{3}} + 4d_{1}^{2} \left(\frac{4}{n_{1}} + \frac{8}{n_{2}} \right) \\ &+ \frac{8}{n_{1}^{2}} + \frac{40}{n_{2}^{2}} + \frac{64}{n_{1}n_{2}} + \frac{256}{n_{1}n_{2}^{2}} + \frac{96}{n_{1}^{2}n_{2}} + \frac{48}{n_{2}^{3}} + \frac{288}{n_{1}n_{2}^{3}} + \frac{352}{n_{1}^{2}n_{2}^{2}} + \frac{384}{n_{1}^{2}n_{2}^{3}} \right) \\ &+ 4d_{1}^{4} \left(\frac{2}{n_{1}} + \frac{8}{n_{2}} + \frac{40}{n_{2}^{2}} + \frac{24}{n_{1}n_{2}} + \frac{48}{n_{2}^{3}} + \frac{88}{n_{2}^{3}} + \frac{88}{n_{1}n_{2}^{2}} + \frac{96}{n_{1}n_{2}^{3}} \right) \right) \right)$$

Next, we look at $\hat{\sigma}_2^2$ From Appendix B, equation (B.5) we have $\hat{\sigma}_{2}^{2} = 2 \sum_{i=1}^{p} \left[\frac{\lambda_{i}^{2}}{\hat{\gamma}_{i}^{2}} - 2 \frac{\lambda_{i}}{\hat{\gamma}_{i}} + 2 \frac{d_{1}^{2} \lambda_{i}}{\hat{\gamma}_{i}^{2}} + 1 \right]$ (C.8) $E(\hat{\sigma}_{2}^{2}) = 2\sum_{i=1}^{p} \left[\frac{E(\hat{\lambda}_{1}^{2})}{E(\hat{\gamma}_{4}^{2})} - 2\frac{E(\hat{\lambda}_{1})}{E(\hat{\gamma}_{4})} + \frac{2d_{1}^{2}E(\hat{\lambda}_{1})}{E(\hat{\gamma}_{4}^{2})} + 1 \right]$ $= 2\sum_{i=1}^{p} \left[(1+\frac{2}{n_1})(1+\frac{2}{n_2})\lambda_i^2 - 2\lambda_i + 1 + 2d_i^2(1+\frac{2}{n_2})\lambda_i \right] (C.9)$ $\left[E\left(\hat{\sigma}_{2}^{2}\right)^{2}\right] = 4E\sum_{i=1}^{p} \left(\frac{\hat{\lambda}_{1}^{2}}{\hat{\gamma}_{1}^{2}} - \frac{2\hat{\lambda}_{1}}{\hat{\gamma}_{1}} + 1 + 2d\frac{\hat{\lambda}_{1}}{\hat{\gamma}_{1}^{2}}\right)^{2} + \frac{\text{cross-product}}{\text{terms}}$ $= 4E\sum_{i=1}^{p} \left(\frac{\hat{\lambda}_{i}^{4}}{\hat{\gamma}_{i}^{4}} + 4\lambda_{i}^{3} \left(\frac{d_{i}^{2}}{\hat{\gamma}_{i}^{4}} - \frac{1}{\hat{\gamma}_{i}^{3}} \right) + 2\hat{\lambda}_{i}^{2} \left(\frac{3}{\hat{\gamma}_{i}^{2}} - \frac{4d_{i}^{2}}{\hat{\gamma}_{i}^{3}} + \frac{2d_{i}^{4}}{\hat{\gamma}_{i}^{4}} \right)$ + $4\hat{\lambda}_{i}\left(\frac{d_{i}^{2}}{\hat{\gamma}_{i}^{2}}-\frac{1}{\hat{\gamma}_{i}}\right)+1\right)$ + cross-product terms $= 4 \sum_{i=1}^{p} \left(\lambda_{i}^{4} \left(1 + \frac{12}{n_{1}} + \frac{44}{n_{1}^{2}} + \frac{48}{n_{1}^{3}} \right) \left(1 + \frac{12}{n_{2}} + \frac{44}{n_{2}^{2}} + \frac{48}{n_{2}^{3}} \right) \right)$ $+ 4\lambda_{1}^{3}\left(1 + \frac{6}{n_{2}} + \frac{8}{n_{2}^{2}}\right)\left(\left(1 + \frac{12}{n_{1}} + \frac{44}{n_{1}^{2}} + \frac{48}{n_{1}^{3}}\right)d_{1}^{2} - \left(1 + \frac{6}{n_{1}} + \frac{8}{n_{1}^{2}}\right)\right)$ $+ 2\lambda_{i}^{2}\left(1 + \frac{2}{n_{2}}\right)\left(3\left(1 + \frac{2}{n_{1}}\right) - 4d_{i}^{2}\left(1 + \frac{6}{n_{1}} + \frac{8}{n_{1}^{2}}\right) + 2d_{i}^{4}\left(1 + \frac{12}{n_{1}}\right)\right)$ $+\frac{44}{n_1^2}+\frac{48}{n_1^3})+4\lambda_i\left(d_i^2\left(1+\frac{2}{n_1}\right)-1\right)+1\right)+\frac{cross-product}{terms}$

$$= 4 \prod_{i=1}^{p} \left[\lambda_{1}^{u} \left(1 + \frac{12}{n_{1}} + \frac{12}{n_{2}} + \frac{144}{n_{1}n_{2}} + \frac{44}{n_{1}^{2}} + \frac{48}{n_{2}^{2}} + \frac{48}{n_{1}^{3}} + \frac{48}{n_{2}^{3}} + \frac{528}{n_{1}^{2}n_{2}} \right] + \frac{123}{n_{1}^{2}n_{2}^{2}} + \frac{123}{n_{1}^{2}n_{2}^{2}} + \frac{1212}{n_{1}^{2}n_{2}^{2}} + \frac{2304}{n_{1}^{2}n_{2}^{3}} \right) + 4\lambda_{1}^{3} \left(\left(1 + \frac{12}{n_{1}} + \frac{6}{n_{2}} + \frac{44}{n_{1}^{2}} + \frac{8}{n_{2}^{2}} + \frac{72}{n_{1}n_{2}} + \frac{2112}{n_{1}^{2}n_{2}^{2}} + \frac{2304}{n_{1}^{2}n_{2}^{3}} \right) \right) + 4\lambda_{1}^{3} \left(\left(1 + \frac{12}{n_{1}} + \frac{6}{n_{2}} + \frac{44}{n_{1}^{2}} + \frac{8}{n_{2}^{2}} + \frac{72}{n_{1}n_{2}} + \frac{264}{n_{1}^{2}n_{2}^{2}} + \frac{96}{n_{2}^{2}n_{1}} + \frac{48}{n_{1}^{3}} + \frac{288}{n_{2}^{2}n_{1}} \right) + 4\lambda_{1}^{3} \left(\left(1 + \frac{12}{n_{1}} + \frac{6}{n_{2}} + \frac{44}{n_{1}^{2}} + \frac{8}{n_{2}^{2}} + \frac{72}{n_{1}n_{2}} + \frac{264}{n_{1}^{2}n_{2}} + \frac{48}{n_{1}^{3}} + \frac{288}{n_{1}^{2}n_{1}} \right) + 4\lambda_{1}^{3} \left(\left(1 + \frac{12}{n_{1}} + \frac{6}{n_{2}} + \frac{44}{n_{1}^{2}} + \frac{8}{n_{2}^{2}} + \frac{72}{n_{1}n_{2}} + \frac{96}{n_{1}^{2}n_{1}} + \frac{48}{n_{1}^{3}} + \frac{288}{n_{1}^{2}n_{1}} \right) + 4\lambda_{1}^{3} \left(\left(1 + \frac{12}{n_{1}} + \frac{6}{n_{1}^{2}} + \frac{44}{n_{1}^{2}} + \frac{72}{n_{1}n_{2}} + \frac{264}{n_{1}^{2}n_{2}} + \frac{96}{n_{1}^{2}n_{1}} + \frac{48}{n_{1}^{2}n_{2}^{2}} + \frac{48}{n_{1}^{2}n_{2}} \right) + 2\lambda_{1}^{2} \left(\left(1 + \frac{6}{n_{1}} + \frac{6}{n_{2}} + \frac{12}{n_{1}^{2}} + \frac{8}{n_{1}^{2}} + \frac{36}{n_{1}^{2}n_{2}} + \frac{48}{n_{1}^{2}n_{1}^{2}} \right) + 2\lambda_{1}^{2} \left(\left(3 + \frac{6}{n_{1}} + \frac{6}{n_{2}} + \frac{12}{n_{1}n_{2}} \right) - 4d_{1}^{2} \left(1 + \frac{2}{n_{1}n_{2}} + \frac{48}{n_{1}^{2}n_{1}^{2}} + \frac{48}{n_{1}^{2}n_{1}^{2}} + \frac{48}{n_{1}^{2}n_{2}} + \frac{48}{n_{1}^{2}n_{1}^{2}} + \frac{48}{n_{1}^{2}n_{1}$$

妆

$$+ 2\lambda_{1}^{2} \left(\left(3 + \frac{2}{n_{1}} + \frac{2}{n_{2}} + \frac{4}{n_{1}n_{2}}\right) + 2d_{1}^{4} \left(1 + \frac{4}{n_{1}} + \frac{4}{n_{1}^{2}}\right) - 4d_{1}^{2} \left(1 + \frac{2}{n_{1}}\right) \right) \right) \\ + 4\lambda_{1} \left(d_{1}^{2} \left(1 + \frac{2}{n_{1}}\right) - 1 \right) + 1 \right] + \frac{\operatorname{cross-product}}{\operatorname{terms}} \quad (C.11) \\ \operatorname{Var}\left(\hat{\sigma}_{2}^{2}\right) = \left[E\left(\hat{\sigma}_{2}^{2}\right)^{2} \right] - \left[E\left(\hat{\sigma}_{2}^{2}\right) \right]^{2} \quad \text{or} \\ \operatorname{Var}\left(\hat{\sigma}_{2}^{2}\right) = 4 \frac{P}{1 + 1} \left[\lambda_{1}^{\mu} \left(\frac{8}{n_{1}} + \frac{8}{n_{2}} + \frac{128}{n_{1}^{2}n_{2}^{2}} + \frac{40}{n_{1}^{2}} + \frac{48}{n_{2}^{2}} + \frac{48}{n_{1}^{3}} + \frac{48}{n_{2}^{3}} + \frac{512}{n_{1}^{2}n_{2}^{2}} \right) \\ + \frac{512}{n_{1}n_{2}^{2}} + \frac{1920}{n_{1}^{2}n_{2}^{2}} + \frac{576}{n_{1}^{3}n_{2}} + \frac{576}{n_{1}n_{2}^{3}} + \frac{2112}{n_{1}^{2}n_{2}^{2}} + \frac{2304}{n_{1}^{2}n_{1}^{2}} \right) + 4\lambda_{1}^{3} \left(d_{1}^{2} \left(\frac{8}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{2}^{2}} + \frac{40}{n_{1}^{2}n_{2}^{2}} + \frac{2112}{n_{1}^{2}n_{1}^{2}} + \frac{2304}{n_{1}^{2}n_{1}^{2}} \right) + 4\lambda_{1}^{3} \left(d_{1}^{2} \left(\frac{8}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{2}^{2}} + \frac{40}{n_{1}^{2}n_{2}^{2}} + \frac{2112}{n_{1}^{2}n_{1}^{2}} + \frac{2304}{n_{1}^{2}n_{1}^{2}} \right) + 4\lambda_{1}^{3} \left(d_{1}^{2} \left(\frac{8}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{2}^{2}} + \frac{40}{n_{1}^{2}n_{2}^{2}} + \frac{2112}{n_{1}^{2}n_{1}^{2}} + \frac{2304}{n_{1}^{2}n_{1}^{2}} \right) + 2\lambda_{1}^{2} \left(d_{1}^{2} \left(\frac{8}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{2}^{2}} + \frac{40}{n_{1}^{2}n_{2}^{2}} + \frac{2112}{n_{1}^{2}n_{1}^{2}} + \frac{2304}{n_{1}^{3}n_{2}^{2}} \right) + 2\lambda_{1}^{2} \left(d_{1}^{2} \left(\frac{8}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{1}^{2}n_{2}^{2}} + \frac{48}{n_{1}^{2}n_{2}^{2}} + \frac{288}{n_{1}^{3}n_{1}^{2}} + \frac{352}{n_{1}^{2}n_{2}^{2}} + \frac{384}{n_{2}^{2}n_{1}^{2}} \right) \right) + 2\lambda_{1}^{2} \left(\left(\frac{4}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{1}^{2}} + \frac{8}{n_{2}^{2}} + \frac{32}{n_{1}^{2}n_{2}^{2}} + \frac{48}{n_{1}^{2}n_{2}^{2}} + \frac{64}{n_{1}^{2}n_{2}^{2}} \right) \right) + 2\lambda_{1}^{2} \left(\left(\frac{4}{n_{1}} + \frac{4}{n_{2}} + \frac{8}{n_{1}^{3}} + \frac{2}{n_{1}^{2}n_{2}} + \frac{2}{n_{1}^{2}n_{2}} + \frac{2}{n_{1}^{3}n_{2}} + \frac{2}{n_{1}^{3}n_{2}^{2}} \right) \right) \right)$$

$$- 4d_{1}^{2} \left(\frac{2}{n_{2}} + \frac{4}{n_{1}} + \frac{12}{n_{1}^{2}n_{2}} + \frac{8}{n_{1}^{2}} + \frac{12}{n_{1}^{2}n$$

Because we do not know the true values of λ_i , we substitute for λ_i in equations (C.7) and (C.12) by $\hat{\lambda}_i$.

Appendix D

Classification Results Tables

<u>8</u>.,

Channels	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper.	Average Algorithm	S.D. of P _{cc} Exper.
1-2	85.4	82.9	81.1	81.0	79.5	82.0	92.4	2.26
1-3	95.9	95.5	92.4	93.5	º3.0	94.0	94.0	1.39
1-4	92.5	95.7	91.6	94.0	93.5	93.4	93.4	1.56
1-5	90.8	96.5	90.7	96.8	94.2	93.8	93.4	3.00
1-6	89.5	96.4	89.3	97.2	93.4	93.2	93.2	3.71
1-7	90.3	96.7	86.5	96.1	94.0	92.7	93.2	4.29
1-8	89.7	96.2	86.4	97.4	95.4	93.0	93.0	4.74
1-9	90.1	95.2	86.5	97.9	95.9	93.1	92.0	4.69
1-10	89.4	96.0	87.5	97.5	95.7	93.2	91.7	4.46
1-11	88.0	95.7	84.1	94.3	96.1	91.6	90.0	5.33
1-12	87.0	94.9	83.7	94.0	94.5	90.8	87.0	5.14

Table D.1 Classification Results of Aircraft, Simulated Data, Using 20 samples per class.

Channels	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper,	Average Algorithm	S.D. of P _{cc} Exper.
1-2	76.3	77.9	84.2	75.1	86.1	79.9	90 - 5	4.92
1-3	93.3	94.7	94.3	91.5	93.1	93.4	90.5	1.25
1-4	92.3	95.2	94.8	95.1	88.5	92.8	89.7	2.67
1-5	92.8	96.3	93.2	95.2	89.8	93.5	89.2	2.52
1-6	90.1	95.0	91.1	94.4	91.4	92.4	88.5	2.16
1-7	91.0	94.0	93.0	89.2	85.0	90.4	87.3	3.56
1-8	88.0	87.6	80.7	83.9	80.9	84.2	86.3	3.51
1-9	89.0	87.4	75.4	88.1	85.5	85.1	83.5	5.56
1-10	74.6	82.8	76.1	79.3	69.9	76.5	79.0	4.87
1-11	67.9	66.1	65.2	69.4	75.0	76.5	73.3	3.87
1-12	70.8	65.1	58.4	69.8	69.8	66.8	68.4	5.18

Table D.2 Classification Results of Aircraft, Simulated Data, Using 13 samples per class.

State of the second state of the

•

Channels	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper.	Average Algorithm	3.D. of P Exper.
1-2	79.0	98.1	97.8	96.1	94.3	93.1	90.4	8.00
1-3	81.5	97.1	86.2	86.9	93.0	89.0	95.0	6.13
1-4	80.0	97.6	81.5	87.9	93.7	98.1	95.0	5.49
1-5	84.6	97.8	80.3	92.9	95.3	90.2	94.8	7.39
1-6	86.4	96.6	77.0	92.8	95.4	89.6	94.2	8.09
1-7	87.9	94.6	78.0	90.2	93.7	88.9	93.7	6.65
1-8	88.8	94.1	80.3	90.7	94.9	89.8	93.3	5.26
1-9	87.2	96.3	80.7	91.2	95.8	90.2	92.6	6.50
1-10	82.6	96.7	80.7	89.0	86.7	87.1	91.6	6.27
1-11	79.5	96.3	79.4	86.8	84.9	85.4	91.0	6.93
1-12	77.1	96.5	78.0	87.0	84.9	84.7	90.0	7.86

Table D.3 (_assification Results of Aircraft, Real Data, Using 20 samples per class.

- . e .

都容赦を記録 せいじゅう うくう

and the statistic sector is a sector of the sector of the

and the second second

Channels	Sample ì	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper.	Average Algorithm	S.D. of P Exper. cc
1-2	95.8	83.4	89.2	78.6	97.6	88.9	90.3	8.06
1-3	90.0	94.5	98.3	83.7	95.5	92.4	90.4	5.70
1-4	90.1	97.6	98.8	84.3	90.8	92.3	89.7	5.94
1-5	91.2	95.0	98.5	84.7	83.3	90.6	88.7	6.52
1-6	93.3	96.5	81.9	87.5	79.6	87.8	87.8	7.21
1-7	94.6	96.1	83.0	88.7	83.8	89.2	87.3	6.01
1-8	83.0	94.7	83.8	87.2	79.9	85.7	86.3	5.65
<u>]</u> 9	64.7	89.4	91.9	85.7	78.8	82.1	83.5	10.91
1-10	62.7	81.6	92.1	83.4	75.6	79.1	80.4	10.90
1-11	62.7	67.9	94.4	79.3	74.8	75.8	76.3	12.12
1-12	66.8	65.1	68.0	70.0	76.1	69.2	71.0	4.25

Table D.4 Classification Results of Aircraft, Real Data, Using 13 samples per class.

Channels	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper.	Average Algorithm	S.D. of P Exper.
1-2	97.4	96.1	97.0	97.9	95.2	96.7	97.8	1.08
1-3	98.7	97.4	98.7	97.9	97.5	98.0	98.1	0.63
1-4	98.7	96.5	98.7	97.9	97.9	97.9	98.1	0.90
1-5	99.0	97.2	98.9	97.6	97.7	98.1	98.0	0.82
1-6	99.0	96.0	98.9	98.1	97.4	97.9	97.8	1.24
1-7	98.9	96.0	98.9	98.2	97.6	97.9	97.3	1.20
1-8	98.7	96.3	98.9	98.2	97.7	98.0	96.9	1.04
1-9	98.7	96.6	98.6	98.2	97.6	97.9	96.5	0.87
1-10	98.7	94.3	98.5	97.1	97.5	97.2	95.6	1.76
1-11	96.0	91.4	98.6	95.5	98.1	95.9	94.5	2.85
1-12	95.6	90.3	98.1	96.1	97.9	95.6	92.1	3.16

Table D.5 Classification Results of Landsat, Multitemporal, Simulated Data, Using 20 samples per class.

Channels	Sample	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper.	Average Algorithm	S.D. of P _{cc} Exper.
				05.9	97.7	97.2	96.5	0.86
1-2	98.0	97.0	97.4	95.0	09 7	98.3	96.5	0.83
1-3	99.0	98.7	96.9	98.2	90.7	07.0	95.6	0.93
1-4	98.0	98.4	96.3	98.0	98.7	97.9	04 7	0.63
1-4	07 4	98.6	97.4	97.9	98.7	98.0	94.1	0.97
1-5	57.4	07.6	96.2	97.5	98.6	97.4	94.0	0.07
1-6	97.1	97.0		98.1	98.7	96.9	92.7	2.00
1-7	97.5	96.7	93.0	,	08.9	96.4	91.8	2.12
1-8	96.0	96.7	93.1	97.2	90.9	05 2	90.8	3.29
1_9	90.6	96.1	93.3	97.2	98.9	93.2	07 0	8.71
1-)	00 /	84.0	82.4	91.3	69.6	83.5	87.0	10.11
1-10	90.4	() 7	77.0	97.1	70.5	75.3	78.0	13.11
1-11	68.1	63.7	,,,,,	077	76.7	72.0	76.7	17.10
1-12	53.4	60.1	72.3	91.1				

Table D.6 Classification Results of Landsat, Multitempoal, Simulated Data, Using 13 samples per class.

n na serie de la construcción de la La construcción de la construcción d La construcción de la construcción d

...

Channels	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper.	Average Algorithm	S.D. of P Exper. cc
1-2	88.3	85.9	91.3	86.0	87.9	87.9	99.0	2.20
1-3	94.7	96.1	95.8	96.6	96.2	95.9	99.2	0.72
1-4	94.3	96.1	96.0	95.8	96.7	95.8	99.1	0.89
1-5	93.5	96.2	95.8	97.0	96.6	95.8	99.3	1.37
1-6	94.0	92.3	95.4	97.0	95.7	94.9	99.1	1.79
1-7	95.4	94.3	93.7	97.1	96.1	95.3	99.0	1.36
1-8	94.0	95.7	93.6	94.8	93.6	94.3	98.8	0.91
1-9	88.8	95.6	94.1	94.0	91.9	92.9	98.4	2.63
1-10	88.3	91.3	93.6	94.3	91.9	91.9	98.2	2.34
1-11	89.2	89.7	93.9	94.5	92.6	92.0	97.0	2.42
1-12	86.7	94.7	93.9	95.8	93.2	92.9	95.5	3.58

Table D 7	
lable D./	Classification Results of Landsat, Multitemporal,
	Real Data, Using 20 samples per class.

We show a product strain and the second strain strai

Channels	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Average Exper.	Average Algorithm	S.D. of P Exper.
1.2	87.5	85.7	91.1	89.7	90.6	88.9	96.3	2.27
1-2	96.0	96.1	96.6	91.8	93.6	94.8	97.0	2.05
1-3	90.0	95.6	97.0	92.3	93.2	94.8	97.6	2.00
1-4	90.1	94 7	95.7	90.7	89.7	93.5	97.0	3.09
1-5	90.0	94.7	90.0	85.4	94.0	91.5	96.3	3.87
1-6	93.1	94.9	94.5	84.4	94.5	92.1	95.4	4.43
1-7	92.2	94.9 05 7	96.1	77.8	94.9	89.5	95.0	7.69
1-8	93.0	07.1	05 3	75.0	93.5	88.8	94.0	8.30
1-9	93.0	07.4	7 7 7	79.5	92.4	87.9	89.6	5.15
1-10	91.8	8/.9	0/./	, , , , , ,			86 7	8.39
1-11	72.1	84.4	79.7	80.1	95.0	82.3	00./	
1-12	50.3	72.2	58.4	74.2	67.1	64.4	78.3	10.00

Table D.8 Classification Results of Landsat, Multitemporal, Real Data, Using 13 samples per class.

The first of a first poor of the strategy terms and the second second second second second second second second

149

.

Appendix E

Computer Program Listings

ORIGINAL PACE IS OF POOR QUALITY

FILE: SWRITE FORTRAN & LARS / PURDUE UNIVERSITY CCCCCC WRITTEN DY. BILL PFAFF EDITED BY: MARWAN MUASHER JUNE 14, 1980 ···· THIS PROGRAM GENERATES SIMULATED DATA BASED DN A CLASSIFICATION MAP OR A CROUND TRUTH MAP EACH PIXEL GENERATED THUS COMES FROM A KNOWN CLASS DISTRIBUTION. THE METHOD USED IS AS FOLLOWS 1. A COUD CLASSIFICATION IS CHOSEN AS A BASE FOR SIMULATED DATA 2. FROM THIS CLASSIFICATION WE KNOW THE NUMBER OF CLASSES, THE CLASS STATISTICS, AND THE CLASS OF EACH PIXEL IN THE AREA CLASSIFIED 3. A STREAM OF UNIFORM RANDOM NUMBERS IS GENERATED FOR EACH CHANNEL. THEY ARE CHANGED IO NORMAL (0,1) DEVIATES. 4. FOR EACH PIXEL, A RANDOM N(0,1) VECTOR IS TRANSFORMED TO DE DISTRIBUTED ACCORDING TO THE CLASS STATISTICS OF THAT PIXEL THIS IS THE SIMULATED DATA VECTOR. 5. AS EACH LINE IS COMPLETED, IT IS WRITTEN TO AN OUTPUT TAPE. TO RUN THE PROGRAM, YOU NEED TO HAVE THE FOLLOWING EXEC FILE ON YOUR DISK: CETDIEK LAPEYR GETDISK LARSYB GETDISK LARSYB GETDISK DVSYS GLOHAL TXTLIB CMSLIB FORTRAN SSP370 FILEDEF 6 PRINTER FILEDEF 16 TERMINAL FILEDEF 12 TAP2 FILEDEF 11 TAP1 (RECFM VS LRECL 1500 BLKSIZE 1500) LOAD SWRITE GLOCOM MMTAPE TAPOP BCDVAL GTSERL GTDATE MFSD RANDU WRTMTX START SWRITE THE PROGRAM WILL ASK FOR INFORMATION SUCH AS TAPE NUMBERS, FILE NUMBERS, ... ETC FROM HERE ON, IT SHOULD BE EASY TO FOLLOW VARIABLES USED IN TPRINT A = COVARIANCE STORAGE FOR FACTORING AREANDEAREA NUMBER OF CLASSIFICATION B = CUVARIANCE STORAGE FOR MULTIPLICATION DATA = DATA FOINT STURAGE DATVALELINE NUMBER AND ROLL PARAMETER ICAL #CALIPHATION INFORMATION IDREC = IDENTIFICATION RECORD STORAGE ISTARTESTARTING POINTS FOR CAUSS LOGDATEDATA POINTS IN LOGICAL FORMAT NOCHANEROMER OF CHAINELS IN CLASSIFICATION NOCHASENUMBER OF CLASSES IN ORIGINAL STATISTICS NOFLDSENUMBER OF TEST FIELDS NOFLSENUMBER OF TODLED CLASSES PNTCLSECLASSIFICATIONS ARRAY Z = STATISTICS STORAGE Ĉ č INTEGER*2 12, INTDAT, ICAL (3), ILIN(2), PNTCLS(1000), ISTAT(4), FETVC3(30) LD3JCAL+1 (1(2), L03DAT(2), LCAL(6), PATDUT(12000) REAL*4 A(70), A2(12), 7(2700), R(12, 12), DATA(12), * RELAH(30, 12), RVAR(30, 12, 12), T2(100), FRGCAL(5, 30) INTEGER*4 JETART(12), ED3, IMEO(17), AMUAND, IDEE(2003), TAPENO, THREE, CLAPHICOC, IMEAN(30, 12), IVAN(30, 12), I2), YES, NO, DATE(3) INTEGER*4 EUG3D, FLGT EQUIVALENCE (FRGCAL(1, 1), IDEE(51)) DATA ED5, 5, AM //ED5 ', 1 0, 0 0 / DATA YES, NO, THREE //YES ', 'ND ', '3'/ \$ \$ £

OFIGINAL PAGE IS OF POOR QUALITY

FILE SHPITE FORTRAN & LARS / PURDUE UNIVERSITY DATA FLGT /'BIN '/ C LUAD TAPES AND READ PARAMETERS Č. ************************* WRITE(16,500) 500 FORMAT(1/5X,'SPECIFY TAPE NUMBER ON WHICH RESULTS FILE IS LOCATED (5X,'(TYPE EIGHT DIGIT TAPE NUMBER)') READ(16,505)INTAP 505 FORMAT(18) WRITE(16,510) 510 FOFMAT(18) (5X, 'SPECIFY FILE NUMBER AT WHICH RESULTS FILE IS LOCATED'/ (5X, '(TYPE THREE DIGIT FILE NUMBER)') READ(16,515)IFILE 515 FOFMAT(13) CALL MMTAPE(INTAP. IFILE.0) WRITE(16,570) 570 FORMAT(7/5X, 'SPECIFY THE TAPE NUMBER ONTO WHICH SIMULATED DATA IS (TO LE WRITTEN'/5X, '(TYPE EIGHT DIGIT TAPE NUMBER)') RLAD(16,575)TAPEND 575 FORMAT(18) WRITE(16,560) WRITE(16,580) 580 FORMAT(5%, 'SPECIFY FILE NUMBER AT WHICH SIMULATED DATA IS TO BE W \$ITEN'/5%, '(TYPE THREE DIGIT FILE NUMBER)') READ(16,585)JFILE 585 FORMAT(13) DITE(14,585)JFILE Set From a file WR THE (16, 590) Set FORMAT (7/5%, 'SPECIFY THE RUN NUMBER FOR THE SIMULATED DATA RUN'/ 1 5%, '(TYPE EIGHT DIGIT RUN NUMBER)') READ(16, 575) RUNNO CALL MOUNT (TAPENO, 12, 'RI') MARG=JFILE=1 IF (MARG LE=0) GD TO 3 DO 3 LIP=1, MARG CALL TOPFF(12) 3 CONTINUE C RCAD(11)I IF(I NE.1) GO TO 310 READ(11)I, J, NOCLAS, NOCHAN, NOFLDS, NOPOOL, (FETVC3(IX), IX=1, NOCHAN) NOCH=((NOCHAN+1)/2)*2 NOC(MP=NOCHAN+(NOCHAN+1)/2 IST(P=NOC(MP=NOCHAN+(NOCHAN+1)/2) IST(P=NOC(MP=NOCHAN+NOPOOL) IEND=ISTOP+NOCHAN*NOPOOL 15 READ(11) I, J, K IF(I LT 3) GD TO 15 IF(K, NE EOS) GO TO 15 RFAD(11) I, J, (Z(IX), IX=1, IEND) DO 17 IX=1, IEND Z2(IX)=Z(IX) 525 FDRFAT(5%, '+DATA SIMULATION USING MECABES EQUATION'', WRITE(6,530) 530 FDPMAT(5%, '++++++++++++++++++++++++++++++') ERITE(6,535) RUNIO, IDREC(3) 535 FORMAT(////5%, 'SIMULATED DATA RUN IS', I9, ' FROM RUN', I9) WRITE(6,537)INFD(4), INFD(5), INFD(7), INFD(8) 537 FORMAT(/5%, 'LINE', I5, ' TO LINE', I5, ' AND COLUMN', I5, ' TO COLUMN', 537 FORMAT(/5%, 'LINE', I5, ' TO LINE', I5, ' AND COLUMN', I5, ' TO COLUMN', 537 FORMAT(/5%, 'LINE', I5, ' TO LINE', I5, ' AND COLUMN', I5, ' TO COLUMN', 537 FORMAT(/5%, 'LINE', I5, ' TO LINE', I5, ' AND COLUMN', I5, ' TO COLUMN', 537 FORMAT(/5), 'LINE', I5, 'TO LINE', I5, 'AND COLORN , 15, (5) ERITE(6, 540) INTA', IFILE
540 FORMAT(/5), 'INPUT RESULTS FILE IS ON TAPE', 19, 'FILE', 14) ERITE(6, 545) TAPELO, FILE
545 FORMAT(/5), 'SIMULATED DATA IS ON TAPE', 19, 'FILE', 14) ERITE(6, 555)
550 FORMAT(/5), 'CHANNELS USED') DU LOO INFINIELS USED') USIDE(6, 555) FETVOR(IX), FROCAL(1, IX), FROCAL(2, IX)
555 FORMAT(5), IP, 2A+F5, 2, '-', F5, 2)
540 CALL GTDATE(DATE) GALL GTDATE(DATE) UNITE(6,060)DATE 565 FORMAT(/5%, 'DATE OF SIMULATION IS ',3A4)

```
FILE: SWRITE
                                            FORTRAN A LARS / PURDUE UNIVERSITY
                    1DEC=1
 Č • • • • •
                                                                                                                             ------
 C FACTOR COVARI
     FACTOR COVARIANCE MATRICES
                                                                                       .....
                                                                                                                                         -----
                   DO 30 IX=1, NOPOOL
IDONE- IBEG+NOCUMP-1
         ID(NE, EIBEG, IDONE
M=0
DD 20 IY=IBEG, IDONE
K=K+1
20 A(K)*2(IY)
CALL MFSD(A, NOCHAN, EPS, IER)
IF(IER EQ -1) GD TD 300
IF(IER GE, 1) GD TD 310
4-0
          17 (12K GE. 17 GO 10

DO 25 IY=IBEG, IDDNE

K=K+1

25 Z(IY)=A(K)

30 IBEG=IBEG+NOCOMP
C GENIERA
C GENIERA
C GENIERA
       GENERATE STARTI' 7 POINTS
                                                                                                                                                                -----
                                                                                *****
       29 WRITE(16,7:)
3: FORMAT(5X, 'DD YOU WANT TO SPECIFIY THE STARTING POINTS FOR THE '/S
4, 'RANDOM NUMBER GENERATOR? (TYPE YES OR NO)')
READ(16,32)INPUT
32 FORMAT(A4)
1F(INPUT EQ.ND) CO TO 36
1F(INPUT EQ.VES) GO TO 33
GO TO 29
33 DO 39 IX=1,NOCHAN
WRITE(16,41)IX
41 FORMAT(5X, 'SPECIFY STARTING POINT FOR CHANNEL', I3/5X, '(TYPE A NIN
4 DIGIT ODD NUMBER)')
READ(16,42)ISTART(IX)
52 FOPMAT(19)
39 CONTINUE
GO TO 43
36 CALL GTSERL(ISERL)
15ERL=(ISERL/10)*0+1
DO 40 I=1,NOCH
15ERL=[SERL+1000000
15TART(1)=ISERL
43 WRITE(6,34)
34 FORMAT(7//5X, 'STARTING POINTS FOR RANDOM NUMBER GENERATOR'//)
DO 44 I=1.NOCHAN
WRITE(6,35)I.ISTART(I)
35 FORMAT(5X, 'STARTING POINT FOR CHANNEL ', I2,' IS ', I9)
44 CONTINUE
                                                                                       C C CASSIFICATIONS
C READ CLASSIFICATIONS
                                                                    ****!
   IDFEC(1)=TAPENO

IDFC(2)=JFILE

IDFC(3)= RUNNO

NULD = IDFEC(5)

IDFEC(5) = NOCHAN

IDFEC(5) = A((NDPNTS + 9)/4)

NO(SAM = IDFEC(6)

IDFC(6) = 4 + ((NDPNTS + 9)/4)

NO(SAM = IDFEC(6)

IDFC(7) = FLGT

DO(141 II=1, 3)

IDFEC(11+16) = DATE(11)

141 CONTINUE

IDFC(20) = NOLINE

DO(145 II=1, NOCHAN

INEW = FETVC3(II)

DO(145 II=1, 5)

FRC(AL(112,11) = FFGCAL(112, INEW)

145 CONTINUE

LIF = FNUCHAN + 1

DO(150 II = LIP.NOLD

DO(150 II = 1,5)

FRGCAL(112,II) = 0.0
                                                                              **********************************
č
00000000000
```

OF POCR QUALITY

FORTRAN A LARS / PURDUE UNIVERSITY FILE: SWRITE C 150 CONTINUE CALL TOPUR(12, 800, IER, IDREC) IF(IER NL 0) URITE(16,234)IER IF(IER NL 0) URITE(16,234)IER IF(IER CT 0) CO TO 310 DD 50 MG 1, NOCLAS CLAPNT(MA 0 DD 50 MG 1, NOCLAS IMEAN(MA, '2'=0 RMEAN(MA, '2'=0 RMEAN(MA, '2'=0 RMEAN(MA, MB, MC)=0 DO 50 MC 1, NOCHAN IVAR(MA, MB, MC)=0 50 RVAR(MA, MB C CENERATE AND WRITE DATA POINTS CONTRACT AND WRITE DATA POINTS CONTRACT AND WRITE DATA POINTS C CONTRACT AND WRITE DATA POINTS 60 I2=1LIN(2) DATOUT(1)=L1(1) DATOUT(2)=L1(2) I2=32767 DATOUT(3)=L1(1) DATOUT(4)=L1(2) DAIDUT(2)=L1(2) IZ=2767 DATUUT(3)=L1(1) DATUUT(4)=L1(2) IZ=0 ICDUNT=ICDUNT+I IZ=PNTCLS(IX) L1(1)=FALSE IPDL=(IZ=1)*NOCHAN DD 65 IZ=1, NOCHAN DD 65 IZ=1, NOCHAN DD 65 IZ=1, NOCHAN DD 65 IZ=1, IY K=K+1 B(IY) IZ)=Z(K) IF(IY)EG IZ) GD TO 65 B(IZ) IY)=0.0 65 CONTINUE DD 70 IY=1, NOCH CALL RANDU(ISTART(IY), NXINP, A2(IY)) ISTAT(IY)=NXINP CALL RANDU(ISTART(IY), NXINP, A(IY)) ISTAT(IY)=NXINP A(IY)=SGRT(=2.*ALOG(A2(IY)))*COS(6.2B318*A(IY)) 70 CONTINUE CLAPHT(IZ)=CLAPHT(IZ)*1 DD 80 IY=1, NOCHAN DATA(IY)=0 IG=NUPDOL NDCOMP+IF(L+IY DD 75 IZ=1, NOCHAN DATA(IY)=DATA(IY)*2(I0) INTDAT=DATA(IY)*2(I0) INTDAT=DATA(IY)

FILE - SWRITE ILE: BWRITE FORTRAN A LARS / PURDUE UNIVERSITY
DD 100 [P*1,NDCLAS
DD 100 [P*1,NDCLAS
DI (CLAPNT(IP) LE 0) GD TD 98
RMLCH(IP,ID)*FLDAT(IMEAN(IP,ID))/FLDAT(CLAPNT(IP))
98 DD 100 1T=ID,NDCHAN
IF(CLAPNT(IP) LE 1) GD TD 100
REPATE(LAT(IVAR(IP,ID))
REVA*FLDAT(ICLAPNT(IP))
REVA*FLDAT(IVAR(IP,ID))
REVA*FLDAT(INTO*IVAR(IP))
REVA*FLDAT(INTO*IVAN, FRGCAL, THREE, FETVC3)
REVA*FLDAT(IVAFINACHAN FRGCAL, THREE, FETVC3)
REVA*FLDAT(IVAFINACHAN REVALATED COVARIANCE MATRIX*)
DD 640 ID=1,NDCHAN
NURC*IIP_INDCHAN
REVAR(IP)
REVAR(IP)
REVAR(IP)
REVAR(IP)
REVARFLOAT(ID)
REVARFLOAT(IP)
REVARFLOAT(ID)
REVARFLOAT(ID)
REVARFLOAT(INCHAN
REVARFLOAT(INCHAN)
REVARFLOATINCE MATRIX*)
REVARFLOATINCE
REVARFLOAT(INCHAN)
RE FORTRAN & LARS / PURDUE UNIVERSITY

 abs
 Format(7775A)
 bitsch(12)
 continued

 b0
 640
 10=1, NOCHAN

 b0
 640
 1N=1, IO

 b0
 640
 1N=1, IO

 b0
 640
 1N=1, IO

 b0
 640
 IN=1, IO

 call
 WRTMTX(A, NOCHAN, FRGCAL, THREE, FETVC3)

 645
 CONTINUE

 call
 TOPEF(12, IER)

 b0
 650

 10
 1X=3, 200

 650
 IDREC(1X)=0

 call
 TOPWR(12, ROO, IER, IDREC)

 IF(IER NE 0)
 WRITE(16, 234) IER

 IF(IER NE 0)
 WRITE(16, 234) IER

 IF(IER NE 0)
 G0 TU 320

 234
 FORMAT(5x, 'ERROR IS', I5)

 č 300 WRITE(6,305) 305 FORMAT(5%,'ERROR -1') 310 WRITE(6,315) 315 FORMAT(5%,'ERROR GT 1') 320 STOP END

CRIGINAL PACE IS OF POOR QUALITY

÷

ORIGINAL PROSENTS

FILE: HUGHES FORTRAN A LARS / PURDUE UNIVERSITY

HUGHES FORTRAN

HUGHES FORTRAN

HUGHES FORTRAN

HUGHES AS INDUT., DECK IN THE READER FILE AS

HUGHES FORMAT

HUGHES AS INDUT., DECK IN THE READER FILE AS

HUGHES FORMAT

- FIIST CARD. NUMBER OF TRAINING SAMPLES OF CLASS 1

(FORMAT I3)

- SIC(HD CARD. NUMBER OF TRAINING SAMPLES OF CLASS 2

(FORMAT I3)

- MARKS AND COVARIANCE MATRICES OF CLASS 1 AND 2 IN

LARS /S FORMAT

THE PROGRAM REQUIRES AS AN OUTPUT THE PROBABILITY OF CORRECT

CLASSIFICATION FOR EACH CHANNEL (FOR CHASNILLI. CHANNEL 1.2

CHARLS 1, 2, 5, ETC.), THE TRANSFORMATION MATRIX AND THE

NEW MEAN AND COVARIANCE MATRICES

THE PROGRAM REQUIRES THE FOLLOWING EXEC FILE

CETDISK IMSL

GLOPAL TATLIB FORTMOD2 CMSLIB DIMSLIB SIMSLIB

LOAD HUGHES

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

NOT NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

NOT NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

NOT NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

NOT NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

NOT NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

NOT NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 1

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMBER OF TRAINING SAMPLES OF CLASS 2

START

LIST OF VARIABLES

NI: NUMB FILE: HUGHES FORTRAN A LARS / PURDUE UNIVERSITY IMPLICIT REAL*8 (A-H.O-Z) REAL*0 SIGNA1(78), SIGMA2(78), AINV(78), WK(192), PS152(12, 12), *WR(160), M1(12), M2(12), PERROR, EGVEC5(12, 1), EGVECT(1, 12), CC(1, 12), *EGVALR(24), EGVECR(PS0), SIGM15(12, 12), AA(1, 1), DEGVEC(12, 12), *EGVAL1(12), HAIACH(12), TEMP1(12), DDE GVC(12, 12), MEANR(2), MEANS(2), *SGMR(2), SGM5(2), GAMAR(2), GAMAS(2), ALPHR(2), ALPHS(2), MEANS(2), *SGMR(2), SGM5(2), GAMAR(2), GAMAS(2), ALPHR(2), ALPHS(2), *SGMR(2), SGM5(2), D(2), DLLTAR(2), DELTAS(2), DIST(2), ERROR(2), *SIGMN(5), SS2NEW(78), ASGMS(2), ASGMA(2), *SIGMN(5(12, 12), DD1(12), DD2(12), TRAN5(12, 12), TRAN51(12, 12), *LAPHLA, WGK(400), MEANS1(12, 2), MEANR1(12, 2), SGMS1(12, 2), SGMR1(12, 2) *DSUC1(12), VSSMA(2) COMPLEXTA EGVAL(12), EGVEC(12, 12), ZN, X1, X2, D1(12), D2(12) EGVIVALENCE (EGVAL(11), EGVALR(1)), (EGVEC(1, 1), EGVECR(1)) *D5 +X1.X2 COUVALENCE (EGVAL(1), EGVALR(1)), (EGVEC(1)) C C C C C C READ NUMBER OF TRAINING SAMPLES OF CLASS C READ NCAN VECTORE OF CLASSES 1 AND 2 C C READ COVARIANCE MATRICES OF CLASS 1 AND C C READ (5, 967)N1 READ NUMBER OF TRAINING SAMPLES OF CLASS 1 AND 2 READ BEAN VECTORE OF CLASSES 1 AND 2 READ COVARIANCE MATRICES OF CLASS 1 AND 2 ***** READ(5,967)N1 READ (5,967)N2 FORMAT(13) READ (5,130)M1 READ (5,130)M2 READ (5,130)SIGMA1 READ (5,130)SIGMA2 FORMAT(2X,5E14.7) N = 12 967 150 = 12 COMPUTE INVERSE OF COVAPIANCE MAIRIX OF CLASS I

Star Press Barris Ha OF FUER CONLETT

FILE: HUGHES

145

148

GO TO 146

5 +(1 0/EGVAL1(1)**4)*(B 0/N1 +B 0/N2 +12B 0/(N1*N2) +40.0/N1**2 6 +40 0/N2**2 +4R 0/N1**3 +1B 0/N2**3 +512 0/(N1**2*N2) 7 +512 0/(N1*N2**2) +1920 0/(N1**2*N2**2) +576 0/(N1**3*N2) 8 +576 0/(N2**3*N1) +2112 0/(N1**2*N2**3) +2112 0/(N1**3*N2**2) 9 +2534 0/(N1**3*N2**3) +4.0*D50R21(1)*(2 + 0/N1 +B 0/N2 + B 0/N1**2 +40 0/N2**2 +64.0/(N1*N2) +256.0/(N1*N2**2) * +96 0/(N1**3*N2) +4E.0/N2**3 +20B 0/(N1*N2**3) +352 0/(N1**2*N2 * +96 0/(N1**2*N2) +4E.0/N2**3 +20B 0/(N1*N2**3) +352 0/(N1**2*N2 * +96 0/(N1**2*N2) +4E.0/N2**3 +20B 0/(N1*N2**3) +352 0/(N1**2*N2 * +96 0/(N1**0**3))) * +96 0/(N1**0**3))) * +96 0/(N1**0**3))) * +96 0/(N1**0**3))) * +1720 0/(N1**0**3) +40 0/N2**2 +40 0/N1**3 +40 0/N2**3 * +96 0/(N1**0**3))) * +1720 0/(N1**0**3) +512 0/(N1**0**3) +1720 0/(N1**2*N2**3) * +122 0/(N1**0**2) +512 0/(N1**2**3) +1720 0/(N1**2*N2**3) * +2112 0/(N1**0*N2**2) +512 0/(N1**2**3) +1720 0/(N1**2*N2**3) * +40 0/N1**3 +256 0/(N2**3*N1) +352 0/(N1**2*N2) +96 0/(N2**2*N1) * +40 0/N1**2 +44 0/(N1*N2) +256 0/(N1**2*N2) +96 0/(N2**2*N1) * +40 0/N1**2 +46 0/(N1**2*N2) +48 0/N1**2*N2) +96 0/(N2**2*N1) * +40 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N2**2*N1) * +40 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N2**2*N1) * +40 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N1**2*N2)) * +64 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N1**2*N2)) * +64 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N1**2*N2)) * +64 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N1**2*N2)) * +16 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N1**2*N2)) * +16 0/(N1**2*N2) +40 0/(N1**2*N2) +48 0/(N1**2*N2) +96 0/(N1**2*N2)) * +16 0/(N1**2*N2) +48 0/N1**3 +88 0/(N1**2*N2) +96 0/(N1**3*N2)) * +16 0/(N1**2*N2) +40 0/N1**3 +88 0/(N1**2*N2) +96 0/(N1**3*N2)) * +16 0/(N1**2*N2) +40 0/N1**3 +88 0/(N1**2*N2) +96 0/(N1**3*N2)) * +16 0/(N1**2*N2) +40 0/N1**3 +88 0/(N1**2*N2) +96 0/(N1**3*N2)) * FORMAT(10X, 'SGMRB = ',F20.4) 974 C CALCULATE MULTIPLICATIVE FACTOR AND NEW $\hat{\tau}$ AND $\hat{\tau}_{c}$ C 979 874 873 PSI=PSI+DLOG(EGVAL1(I))+((DD2(I)-DD1(I))++2)/(EGVAL1(I)-1.0) PSI=PSI+DLOG(EGVAL1(I))+((DD2(I)-DD1(I))**2)/(EGVAL1(I)) DIST(J)=PSI-(CR(J)-CS(J)) CONTINUE DD 145 K =1,2 IF (DIST(K),LT 0.0)GO TO 147 IF (DELTAR(K),LT 0.0)GO TO 147 IF (DELTAR(K))* SD TO 146 ERROR(K)=1.0 GO TO 146

157

FORTRAN A LARS / PURDUE UNIVERSITY

ORIGINAL PAGE IS OF POOR QUALITY

```
FILE HUGHES FORTRAN A LARS / PURDUE UNIN REITY
С
             CALL LINV20(SIGMA1, N, AINV, IDGT, EE1, EE2, WK, IER)
WRITE(6,117)IER
FORMAT(7 7, I3)
  117
C**
            CONPUTE IN A RSE OF COVARIANCE MATRIX 1 MULTIPLIED
BY COVARIANCE MATRIX 2
č
            CALL VMULSS (AINV, SIGMAZ, N, PS152, N)
C
                                                                                     ***************
COMPUTE EIG(NVALUES AND EIGENVECTORS OF (INVERSE(
SIGMA1)) (SIGMA2)
            CALL_EIGRF(PS1S2, N, N, 2, EGVALR, EGVECR, N, WR, IERR)
WRITE(6,117)IERR
WRITE(6,126)WR(1)
FORMAT('', F6.1)
126
C***
C
C***
                                                                                                     ********
         NORMALIZING EIGENVECTORS (SEE FUKUNACA,
PAGE 35)
                                                                                     ************
           CALL VCVISF (SIGMA1, N, SIGMIS, N)

CALL VCVISF (SIGMA2, N, SIGMPS, N)

D0 10 I = 1, N

D0 70 J = 1, N

EGVECT(I, J) = DRFAL(EGVEC(J, I))

EGVECT(I, J) = DRFAL(EGVEC(J, I))

CONTINUE

M = N

CALL VMULFF(EGVECT, SIGMIS, 1, M, NN, 1, N, CC, 1, IEER)

WRITE(6, 126) IEER

CALL VMULFF(CC, EGVECD, 1, M, NN, 1, N, AA, 1, IIER)

WRITE(6, 126) IER

CALL VMULFF(CC, EGVECD, 1, M, NN, 1, N, AA, 1, IIER)

WRITE(6, 126) IER

CALL VMULFF(CC, EGVECD, 1, M, NN, 1, N, AA, 1, IIER)

WRITE(6, 126) IER

CALL VMULFF(CC, EGVECD, 1, M, NN, 1, N, AA, 1, IIER)

WRITE(6, 126) IER

CALL VMULFF(CC, EGVECD, 1, M, NN, 1, N, AA, 1, IIER)

WRITE(6, 126) IER

CONTINUE
č
  20
  30
    10
              CONTINUE
C***
                           **********
C
C
C CALCULATE
C CALCULATE
          CALCULATE NEW MEAN VECTOR DI = EGVEC+MI
            90
            CONTINUE
****************
           D0 95 I =1,N

D0 95 J =1,N

DEGVEC(1,J) = DREAL(EGVEC(J, 1))

EGV/L(1) = DREAL(EGVAL(I))

DDEGVE(I,J) = DREAL(EGVEC(I,J))

D1(1)=EGVEC(J,I)*M1(J)+D1(I)

DA(1)=EGVEC(J,I)*M1(J)+D1(I)

DA(1)=EGVEC(J,I)*M1(J)+D2(I)

TKANG1(I,J)=0

CUMINUE

FLOMAT(SF14,7)

D0 T77 I =1,N

DD1(I)=DETAL(D1(I))

D1(I)=DETAL(D1(I))

CLOTINUE
  95
183
 777
              CLISTINUE
č***
         *****************
```

```
FILE: HUGHES FORTRAN & LARS / PURDUE UNIVERSITY
 ссс**
          DRDER THE EIGENVALUES AND EIGENVECTORS ACCORDING TO MAXIMUM EIGENVALUE
                                                          *******************
  D0 120 1=1.N

D0 120 J=1.N

IF (EGVAL1(I)-EGVAL1(J))120.120.131

131 TEMP=ED1(I)

TEMP=DD2(I)

EGVAL1(I)=EGVAL1(J)

DD1(J)=D1(J)

DD2(I)=DD2(J)

EGVAL1(J)=TEMP

DD1(J)=TEMP

DD1(J)=TEMP

DD132 K=1.N

TEMP1(K)=DDEGVC(K, J)

DDEGVC(K, J)=TEMP1(K)

132 CONTINUE

120 CONTINUE
120 CONTINUE
CONTINUE
C INITIALIZE ALL PARAMETERS UNDER CONSIDERATION
C INITIALIZE ALL PARAMETERS UNDER CONSIDERATION
              WRITE(6, 136)

DD 134 J=1,N

DD 134 J=1,N

TRANS(I,J)=DDEGVC(J,I)

CONTINUE

DD 135 II =1,2

MEANS(II)=0.0

SGMS(II)=0.0

GAMAS(II)=0.0

GAMAS(II)=0.0

ALPHR(II)=0.0

ALPHR(II)=0.0

DELTAS(II)=0.0

DELTAS(II)=0.0

CR(II)=0.0

CS(II)=0.0

CS(II)=0.0

CS(II)=0.0

CS(II)=0.0

CONTINUE
   134
               CONTINUE
   135
FCRMAT(' ',10X, 'FIRST N DIMENSIONS',10X, 'PROBABILITY DF ERROR')
DD 140 I = 1,N
A(1)=1 0-1 0/EGVAL1(I)
B(1)=(DD1(I)-DD2(I))/(EGVAL1(I)-1.0)
A(2)=EGVAL1(I)-1
0
B(2)=(DSGRT(EGVAL1(I))*(DD1(I)-DD2(I)))/(EGVAL1(I)-1.0)
DSGR21(I)=(DD1(I)-DD2(I))**2
   136
C CALCULATE VAR ( ( ) AND VAR ( ) )
                                                                                                     ******
                                                                                                                        *****
 Ĉ
            VSGMA(1)=VSGMA(1)+4.0+((2.0/EGVAL1(1)++2)+(4.0/N1 +4.0/N2
1.+8.0/(N1+N2)) -(4.0/EGVAL1(1)++3) + (4.0/N1 +4.0/N2 +
2.8.0/(N1++2)+8.0/(N2++2+32.0/(N1+N2)) + (4.0/(N1+N2+2))
3.+.48.0/(N1++2+N2) +54.0/(N1+N2) +4.0/(N1+N2+2))
4.(1.0/N1 +2.0/N2 +6.0/(N1+N2) +4.0/(N2++2.+8.0/(N1+N2++2)))
```

159

ORIGINAL PACE IS OF POOR QUALITY

FILE HUCHES FORTRAN A LARS / PURDUE UNIVERSITY 147 190 CONTINUE C PRINT TRANSFORMATION MATRIX AND NEW MEAN AND COVARIANCE C MATRICES C HOLTE(4 BIB) WRITE(6,919) FORMAT(10X, 'TRANSFORMATION VECTOR') WRITE(6,103)((TRANS(1,J),J,1,N),I=1,N) WRITE(6,103)((TRANS(1,J),J,1,N),I=1,N) FORMAT(10X, 'NEW MEAN VECTORS AND COVRIANCE MATRICES OF CLASS 1 AND 2') WRITE(6,165)(DD1(I),I=1,N) WRITE(6,165)(DD2(I),I=1,N) FORMAT('MN',5E14.7) DO 746 J=1,N IF(I NE J)GD TO 747 SSINEW(I+(1*(I-1))/2)=1.0 SS2NEW(I+(1*(I-1))/2)=EGVAL1(I) GO TO 746 SSINEW(I+(J*(J-1))/2)=0.0 SSNEW(I+(J*(J-1))/2)=0.0 CONTINUE NMN=N*(N+1)/2 WPITE(6,175)(SSINEW(I),I=1,NMN) WRITE(6,175)(SS2NEW(I),I=1,NMN) FORMAT('CV',5E14.7) STOP END 919 920 921 145 747 746 748 175 432 STOP

•

ORIGINAL PROD

Appendix F

. . Description of Data Sets For Experiments

OF POOL OUALIS

F.1 Training and Test Fields for Aircraft, Simulated

Data Set (Tape 203, file 3)

Training Fields

CLASS CORN RUN(71053900), LINE(304, 326, 2), COL(109, 133, 2) RUN(71053900), LINE(512, 528, 1), COL(87, 93, 1) RUN(71053900), LINE(620, 636, 1), COL(107, 123, 2) RUN(71053900), LINE(656, 676, 2), COL(33, 59, 2) CLASS FOREST RUN(71053900), LINE(798, 812, 1), COL(141, 161, 2) RUN(71053900), LINE(704, 720, 1), COL(147, 155, 1) RUN(71053900), LINE(726, 736, 1), COL(81, 95, 1)

Test Fields (Also Area Classified)

TEST CORN
RUN(71053900), LINE(143, 154, 1), COL(42, 57, 1)
RUN(71053900), LINE(305, 318, 1), COL(116, 132, 1)
RUN(71053900), LINE(403, 413, 1), COL(17, 33, 1)
RUN(71053900), LINE(643, 657, 1), COL(121, 127, 1)
RUN (71053900), LINE (684, 691, 1), COL (11, 30, 1)
RUN (71053900), LINE (857, 866, 1), COL (34, 53, 1)
TEST EDREST
PUN(71053900) + INE(424, 430, 1), COU(141, 173, 1)
DUN(71053000) LINE(E7 53) 1) COL(181) 1/3, 1/
RUN(71053700), LINE(521, 531, 1), CUL(142, 162, 1)
RUN(/1053400), LINE(/11, /28, 1), CUL(144, 158, 1)
RUN(71053900), LINE(769, 779, 1), COL(127, 148, 1)
RUN(71053900), LINE(837,851,1), COL(155,162,1)
RUN(71053900), LINE(923,931,1), COL(70,79,1)

OF PG-32 (1200) - 163

F.2 Training and Test Fields for Aircraft, Real

Data Set (Tape 203, file 1)

Training Fields

CLASS CORN RUN(71053900), LINE(304, 326, 2), COL(109, 133, 2) RUN(71053900), LINE(512, 528, 1), COL(87, 93, 1) RUN(71053900), LINE(620, 636, 1), COL(107, 123, 2) RUN(71053900), LINE(656, 676, 2), COL(33, 59, 2) CLASS FOREST RUN(71053900), LINE(798, 812, 1), COL(141, 161, 2) RUN(71053900), LINE(704, 720, 1), COL(147, 155, 1) RUN(71053900), LINE(726, 736, 1), COL(81, 95, 1)

Test Fields (Also Area Classified)

TEST CORN RUN(71053900), LINE(227, 247, 1), COL(81, 96, 1) RUN(71053900), LINE(334, 351, 1), COL(66, 100, 3) RUN(71053900), LINE(452, 474, 2), COL(108, 119, 1) RUN(71053900), LINE(597, 611, 1), COL(137, 153, 2) RUN(71053900), LINE(597, 611, 1), COL(101, 128, 2) RUN(71053900), LINE(646, 664, 1), COL(101, 128, 2) RUN(71053900), LINE(711, 721, 1), COL(102, 113, 1) TEST FOREST RUN(71053900), LINE(241, 249, 1), COL(102, 113, 1) TUN(71053900), LINE(509, 527, 1), COL(181, 193, 1) RUN(71053900), LINE(729, 751, 2), COL(201, 217, 1) RUN(71053900), LINE(76, 803, 2), COL(191, 203, 2) RUN(71053900), LINE(833, 855, 2), COL(191, 171, 2) RUN(71053900), LINE(989, 1005, 1), COL(141, 155, 2)

ORIGINAL BARRIER

÷ ,

F.3 Training and Test Fields for Landsat, Multitemporal,

Simulated Data Set (Tape 203, file 6)

Training Fields

CLASS CORN 78843016 78843016 78843016 78843016 CLASS COVP	25 62 30 91	32 67 33 97	1 1 1 1	33 133 87 79	42 141 102 86	1 1 1 1
78843016	9	12	1	61	77	1
78843016	74	82	1	51	64	1
78843016	110	117	1	167	172	1

Test Fields (Also Area Classified)

TEST CORN
RUN(78843016), LINE(2,12,1), COL(30,34,1)
RUN (78843016), LINE (38, 46, 1), COL (18, 26, 1)
RUN (78843016), LINE (55, 58, 1), COL (103, 117, 1)
RUN (78843016), LINE (16, 22, 1), COL (123, 127, 1)
RUN (78843016), LINE (70, 73, 1), COL (80, 89, 1)
RUN (78943014), LINE (85, 93, 1), COL (47, 50, 1)
RUN (78843016), LINE (102, 104, 1), COL (140, 155, 1)
RUN(78843014), LINE(107, 115, 1), COL(11, 15, 1)
1631 30106AN3 BUNU78042046A LINE(1 4 1) COL(01 100 1)
RUN(78843010), LINE(16, 20, 1), CUL(36, 70, 1)
RUN(78843016), LINE(32,34,1), CUL(114,126,1)
RUN(/8843016), LINE(49, 51, 1), CDL(113, 125, 1)
RUN(78843016), LINE(76, 84, 1), COL(31, 40, 1)
RUN(78843016), LINE(99, 106, 1), COL(127, 132, 1)
RUN(78843016), LINE(106, 114, 1), COL(53, 59, 1)

ORIGINAL PAGE 19 OF FOOR QUALITY

F.4 Training and Test Fields for Landsat, Multitemporal, Real Data Set (Tape 203, file 5)

Training Fields

CLASS CORN 78843016 78843016 78843016 78843016 CLASS SOVE	26 91 62 30	32 98 67 34	1 1 1	32 79 134 91	42 86 141 102	1 1 1 1
78843016	9	13	1	68	78	1
78843016	74	82	1	51	63	1
78843016	100	105	1	120	132	1

Test Fields (Also Area Classified)

TEST CORN
RUN(78843016), LINE(2,11,1), COL(27,32,1)
RUN(78843016), LINE(38,46,1), COL(19,25,1)
RUN(78843016),LINE(103,106,1),COL(140,156,1)
RUN(78843016), LINE(101, 115, 1), COL(12, 17, 1)
RUN(78843016), LINE(78,86,1), COL(124,128,1)
RUN (78843016), LINE (67, 74, 1), COL (94, 98, 1)
RUN (78843016), LINE (35, 41, 1), COL (123, 127, 1)
TEST SOYBEANS
RUN (78843016), LINE (41, 44, 1), COL (67, 79, 1)
RUN (78843016), LINE (79, 84, 1), COL (31, 40, 1)
RUN (78843016), LINE (106, 114, 1), COL (54, 59, 1)
RUN (78843016), LINE (44, 51, 1), COL (118, 123, 1)
RUN (78843016), LINE (1, 4, 1), COL (90, 100, 1)
RUN (78843016), LINE (109, 113, 1), COL (132, 147, 1)
RUN(78843016), LINE(44, 47, 1), COL(155, 161, 1)

- v

ORIGINAL PAGE IS OF POOR QUALITY

F.5 Training and Test Fields for Aircraft Binary Tree

Example (Tape 203, file 1)

Training Fields

CLASS WHT1 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 CLASS WHT2	111440034940017 10017	62789 62289 66229 66229 66229 6623 6623 6623 663 663 663 663 663 663	6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 6627 895 7895 7895 7895 7895 7895 7895 7895		164 164 159 1637 775 167 159 1633 1633	1649 1667 1667 775 1671 1561 1633	
71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	34780781227071 12227071	3146 3117 31199457 32257 324599 446691	314 317 317 319 322 3228 3228 3228 3228 3228 3228 3228	1 1 1 1 1 1 1 1 1 1	1669777757897775	16369777757 1657777577577577577577577577577577557	
71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	2137458900369	480 8823 8887 8887 8890 8975 8975 8975 8994	484 880 882 883 885 885 887 887 887 887 887 887 887 887	111111111111111111111111111111111111111	5244834582 1224834582 1224834582 12333282 12334435	55 1326 1268 1288 1334 1358 1398 1395 1282 1394 1395 1395 1395	

ORIGINAL PACT IS OF POCR QUALITY		167					
CLASS PAS1 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	224 115670025	441822245 1012245 101245 100145 100177 10012	402 417 1012 1012 1012 1014 1015 1017 1017 1018 1020		15739 1007 1007 10032 11527 1107	157 1539 101 1027 101 1022 1013 1023 1152 1152 107	
71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	000000000000000000000000000000000000000	43889999355666575555555555555555555555555555	418 589 589 589 593 595 596 596 596 597 597	111111111111111111111111111111111111111	147 665 679 711 671 597 673 673	147 65 67 75 67 57 67 57 67 57 63	งสงกงกงกงกงกงกง
CLASS SDY 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	432-1220 200	424 3358 4880 312 312	424 3352 488 500 312 312	NNNNNNN	125 145 123 1237 47	125 165 123 1237 637 67	20000000000000000000000000000000000000
71053900 71053900 71053900 71053900 71053900 71053900	57 11 41 23	424 426 426 502	424 426 426 502	งกงกง	131 113 137 137 137	131 113 137 137 119	2005- 2005- 2005- 2005- 2005- 2005-
71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	8017115339986449	552225224 552225224 552224 552224 552224 54 54 54 54 54 54 54 54 54 54 54 54 54	51213562267 55225562267 66652 552227 66662 66672 66682 66672 66682 66672 66682 66672 66682 66672 66682 66672 66682 66672 66672 66672 66672 66726 67726 6776767676 677676 67767676 67767676 67767676 6776767676 6776767676767676767676767676767676767676		98731 98731 12339 119055551 11935551	937 921 1233 119 123 119 355 41 355 41	
71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	1370347050018 0310310	731 709 711 726 722 735 805 805 805 805 718	731 709 711 726 726 735 803 805 809 709 718		85417055275 149955275 14441 151	854 151 147 955 849 1495 1495 1495 1495 1495 1495 1495 14	
ONIC THE PURCHARTY

<u></u>	SS WAT 53700 53700 53700 53700 53700 53700 53700 53700 53700 53700 53700 53700	589133448011 11	8812 8892 9933 9933 9933 9933 9943 9944 9944 99	R9726889793391 997388979733443	111111111111111111111111111111111111111	1624 164 134 143 144 140 130	165 164 137 143 143 1440 138 140	
71053900 11 944 944 1 140 140 1N 71053900 14 947 947 1 141 141 1N 71053900 15 949 947 1 141 141 1N	53900 53900 53900	14	944 947 940	944 947 940	1	140	140	1N5- 1N5-

Test Fields (Alse Area Classified)

۰.

TEST 7105	WHEAT		304	312	1	155	161	1	WHEATCUT
7105	3700	UU6 U6	839 854	848 861	1	73	70	1	WHEATCUT
7105	3700 3700	UU7 HH3	829 619	851 641	22	73 151	91 161	2	WHEAT
7105	3900	<u>GG2</u>	569	575	1	145 81	148	1	DATSCUT
TEST	HAY	770			-			•	
7105	3900	L8	879	923	12	17 85	9 9	2	HAY
7105	3700	C4 62	252	275	2	33	35	1	HAY
7105	ŽŹÓŎŎ	Č Š	713	715	ī	37	50	ī	HAY
7105	3700	BBF	313	327	1	173	185	1	HAY
TEST 7105	9ASTU	L2	589	599	1	77	93	1	PASTURE
7105	3700	Z21	1021	1031	1	103	117	1	PASTURE
<u>7105</u>	<u>ăźóŏ</u>	ĬŻ	669	675	ĩ	<u>ĭ</u> ġı	123	2	PASTURE
7105	3700	ння	683	693	1	97 97	129	2	PASTURE
7105	3700	EE5 Z20	421 423	439 445	22	177	191 27	1	PASTURE
TEST	SOYBE	ANS	500	440	-		1.77	-	COVPEANE
7105	3900	G4	649	687	2	77	83	1	SOYBEANS
7105	3900	RR2 115	861 649	867 671	12	123	149 191	2	SOYBEANS
7105	3900	ōġž	479	519	2	105	129	Ž	SOYBEANS
7105	3900	Z9	205	231	2	195	211	ž	SOYBEANS

CHE POUR QUALITY

,

-

TEST CORN 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	A3 A5 C1 F5 DD3 HH1 JJ1 Z15 F6 Z16	227 2283 374 352 711 481 373 305	247 2495 3497 4511 515 387 327	1111211212	81 49 67 108 137 102 3 47 191	9695991153 1153 1112795		CORN CORN CORN CORN CORN CORN CORN CORN
71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	A10 Z6 Z3 EE4 RR4 HH10 M3 Z18	241 729 765 522 833 765 783 375	247 751 8255 8555 795 387	12211211	27 201 191 155 151 139 49 191	45 217 203 159 171 159 81 201	NNNN	FOREST FOREST FOREST FOREST FOREST FOREST FOREST
71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900 71053900	49 47 47 47 47 47 47 47 47 47 45 44	205 817 221 1000 1010 949 849 873 977 1041	209 819 224 1004 1014 973 855 879 983 1047	111111111111111111111111111111111111111	34 49 51 36 20 113 113 11	38 529 539 135 1209 1209 119 15	111111111111111111111111111111111111111	PONDWATR PONDWATR PONDWATR WATER WATER WATER WATER WATER WATER WATER WATER

119

÷

7,7

F.6 Training and Test Fields for Landsat, Multitemporal

Binary Tree Example	(Tape	203,	file	5)
---------------------	-------	------	------	----

Training Fields

CLASS CORN 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 CLASS SDYBEANS	000000000000000000000000000000000000000	000000004400004 0000000000000000000000	870000004450004 000000004450004		335724459447 13374359447 133713740 10902 1000 1000	3572459447 13714942 19942	1 1 1 1 1 1 1 1 1 1 1 1 1 1
78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016	0000000000000	11 13 74 75 76 76 780 81 800 101	11 13 74 75 76 76 77 81 82 100 101	111111111111111111111111111111111111111	67273261309850 1111	6727 54556 54509 855 120 120	
78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016 78843016	000000000000000000000000000000000000000	52223511245217	555555999995	111111111111111111111111111111111111111	154 150 158 160 158 160 177 88 90 90 90 90	154 154 160 158 161 180 182 177 178 188 39 50 49	

OF 1 201. ALL MEY

-

Test Fields (Also Area Classified)

TEST_CORN RUN(78843016),LINE(2,11,1),COL(27,32,1) RUN(78843016),LINE(38,46,1),COL(19,25,1) RUN(78843016),LINE(103,106,1),COL(140,156,1) RUN(78843016),LINE(101,115,1),COL(124,128,1) RUN(78843016),LINE(78,86,1),COL(124,128,1) RUN(78843016),LINE(67,74,1),COL(94,98,1) RUN(78843016),LINE(35,41,1),COL(123,127,1) TEST_SOYBEANS RUN(78843016),LINE(41,44,1),COL(67,79,1) RUN(78843016),LINE(106,114,1),COL(31,40,1) RUN(78843016),LINE(106,114,1),COL(116,123,1) RUN(78843016),LINE(1,4,1),COL(116,123,1) RUN(78843016),LINE(1,4,1),COL(112,147,1) RUN(78843016),LINE(109,113,1),COL(135,161,1) TEST_ELSE RUN(78843016),LINE(33,42,1),COL(137,141,1) RUN(78843016),LINE(54,57,1),COL(137,141,1) RUN(78843016),LINE(54,57,1),COL(136,149,1) RUN(78843016),LINE(55,59,1),COL(136,149,1) RUN(78843016),LINE(108,114,1),COL(83,89,1)

į

1

171